

DIVERGENT COMPUTE • AI ECONOMIC THINK TANK

First Principles.

How AI actually works — from a token to the \$500-billion question.

I • Foundations

II • Models

III • Inference & systems

IV • Building with AI

V • The frontier & the industry

VI • Best practices & tools

34 chapters • six parts • every code block run before publishing
v1.0 • CC-BY 4.0 • divergentcompute.com/first-principles

Contents

I • FOUNDATIONS

- 01 What is a token?
- 02 What is an embedding?
- 03 What is a neural network?
- 04 What is attention?
- 05 What is a transformer?
- 06 Parameters & weights
- 07 The context window

II • MODELS

- 08 What is an LLM?
- 09 How models are pretrained
- 10 Fine-tuning & RLHF
- 11 Differences between LLMs
- 12 Multimodal models
- 13 Quantization & distillation

III • INFERENCE & SYSTEMS

- 14 What is inference?
- 15 Latency, throughput, tokens/sec
- 16 Why AI needs GPUs
- 17 Memory, HBM & the memory wall
- 18 KV cache & batching
- 19 The data-center cluster

IV • BUILDING WITH AI

- 20 Prompting
- 21 What is RAG?
- 22 Agents & tool use
- 23 RAG vs fine-tune vs prompt
- 24 Evals
- 25 Vector search

V • THE FRONTIER & THE INDUSTRY

- 26 Scaling laws

- 27 The foundation-model labs
- 28 The compute supply chain
- 29 The economics of a token
- 30 Multi-agent & what comes next

VI • BEST PRACTICES & TOOLS

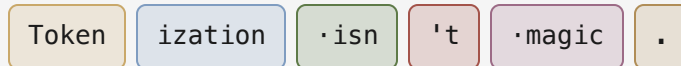
- 31 Choosing a model
- 32 Cost optimization
- 33 Safety, evals & guardrails
- 34 The tool landscape

What is a token?

A token is the basic unit of text an AI model reads and writes — **not a word, not a letter, but a chunk**. Models don't see language; they see tokens.

Read at your depth: [01 The answer](#) · [02 Intuition](#) · [03 Mechanics](#) · [04 The math](#) · [05 The code](#) · [06 The economics](#) · [07 Sources](#) — the deep layers expand on click.

Tokenization isn't magic.



4 words become 6 tokens. The model reads the chunks on the bottom, never the sentence on top.

The model never sees your sentence

Before a language model reads a single thing, your text is chopped into **tokens** by a piece of software called a *tokenizer*. A token can be a whole word (`cat`), a fragment of one (`token` + `ization`), a single character, a space, or a punctuation mark. Common words usually become one token; rare words get split into pieces.

The model then works entirely in tokens — it predicts the next token, over and over, and the tokens are stitched back into text at the end. Three things you'll meet everywhere downstream are all measured in tokens: **the context window** (how much it can read at once), **the price** (APIs bill per token), and **the speed** (tokens per second).

A useful rule of thumb for English: **1 token \approx 4 characters \approx 0.75 words**. So a 1,000-word essay is roughly 1,300 tokens.

TRY IT — THE REAL GPT TOKENIZER, LIVE

Tokenization is how a model reads text. Try "antidisestablishmentarianism".

Token ization ·is ·how ·a ·model ·reads ·text · ·Try ·"
ant idis establish ment arian ism ".

18 **75** **4.2** **\$0.000045**
TOKENS CHARACTERS CHARS / TOKEN COST AT \$2.50/1M

Real GPT-4 tokenizer (cl100k_base), running in your browser.

How the chunks get made: Byte-Pair Encoding

Almost every modern model tokenizes with a variant of **Byte-Pair Encoding (BPE)**. The idea is borrowed from a 1994 data-compression trick and is delightfully simple. You build the vocabulary *by merging*:

- **Start** with the smallest possible units — individual bytes (so any character in any language is representable).
- **Count** every adjacent pair of units across a huge text corpus.
- **Merge** the single most frequent pair into one new unit, and add it to the vocabulary.
- **Repeat** thousands of times until the vocabulary hits a target size — GPT-2 stopped at ~50,000, GPT-4's `cl100k` at ~100,000, the newest at ~200,000.

The result is a learned list of merges. To tokenize new text, the tokenizer greedily applies those merges: frequent strings like `the` survive as one token, while a rare word like `antidisestablishmentarianism` gets rebuilt from several pieces. Each final token maps to an integer ID, and that ID indexes a row in the model's [embedding table](#) — which is the next chapter. Text → tokens → IDs → vectors.

Two consequences worth holding onto: spaces are usually glued to the *front* of the next word (`token` , with a leading space, is a different token from `token`), and the same word can tokenize differently depending on what's around it. That is why token counts feel slightly unpredictable — and why you should measure, not guess.

Vocabulary, merges, and the cost it sets

Let the vocabulary be a finite set of size $|V|$. A tokenizer is a function that maps a string s to a sequence of token IDs $t = (t_1, \dots, t_n)$, $t_i \in \{0, \dots, |V| - 1\}$, where n is the **token length** of s .

BPE builds V greedily. Beginning from the base alphabet V_0 (the 256 bytes), it repeatedly chooses the adjacent pair with the highest corpus frequency and adds its concatenation to the vocabulary:

$$(a^*, b^*) = \arg \max_{(a,b)} \text{freq}(a, b), \quad V_{k+1} = V_k \cup \{a^*b^*\}$$

stopping when $|V_k|$ reaches the target. There is no clean global optimum — BPE is a greedy compressor — but it reliably shortens sequences for a given vocabulary budget.

That budget is the whole game, because tokenization sets two of the model's biggest costs:

- **Sequence length drives compute.** Self-attention compares every token with every other token, so its cost per layer grows as $O(n^2 \cdot d)$ in the number of tokens n (with model width d). Halving your token count roughly *quarters* the attention work.
- **Vocabulary size drives parameters.** The embedding and output layers each hold $|V| \cdot d$ parameters. A bigger vocabulary makes sequences shorter (good) but those two matrices larger (costly).

So a tokenizer is a dial between n and $|V|$ — between sequence length and parameter count — and every model picks a point on it. Hold the $O(n^2)$; it returns in

layer 06.

04 THE CODE

BPE in a few lines — and the real thing

A minimal BPE trainer. It learns merges from a toy corpus exactly as described above — runnable as-is in Python.

`train_bpe.py`

```

from collections import Counter

def get_pairs(tokens):
    return Counter(zip(tokens, tokens[1:]))

def train_bpe(words, num_merges):
    # words: list of strings; start from characters
    corpus = [list(w) for w in words]
    merges = []
    for _ in range(num_merges):
        pairs = Counter()
        for tok in corpus:
            pairs.update(get_pairs(tok))
        if not pairs:
            break
        (a, b), _ = pairs.most_common(1)[0] # most frequent adjacent pair
        merges.append((a, b))
        # merge that pair everywhere
        corpus = [merge(tok, a, b) for tok in corpus]
    return merges

def merge(tok, a, b):
    out, i = [], 0
    while i < len(tok):
        if i < len(tok) - 1 and tok[i] == a and tok[i+1] == b:
            out.append(a + b); i += 2
        else:
            out.append(tok[i]); i += 1
    return out

words = ["token"] * 6 + ["tokenizer"] * 3 + ["tokenization"] * 2
print(train_bpe(words, num_merges=5))
# -> [('t','o'), ('to','k'), ('tok','e'), ('toke','n'), ('token','i')] (piece

```

In production you don't train your own — you call the real one. OpenAI's `tiktoken` gives the exact tokenizer the models use:

`count_tokens.py`

```
import tiktoken

enc = tiktoken.get_encoding("cl100k_base")    # GPT-4 / GPT-3.5
ids = enc.encode("Tokenization isn't magic.")
print(ids)                                   # [3404, 2065, 4536, 956, 11204, 13]
print(len(ids))                             # 6 tokens
print([enc.decode([i]) for i in ids])
# ['Token', 'ization', ' isn', "'t", ' magic', '.']
```

The live widget at the top runs exactly this tokenizer in your browser — type into it and watch the chunks change.

The token is the atom of the AI economy

MECHANICS → MONEY

Everything you just read has a price tag. The token isn't only the unit of *text* — it's the unit of **cost**. Models are billed per token; a frontier model runs roughly **2–5 per million input tokens**. One million tokens is roughly 750,000 words — several thousand pages. Multiply by hundreds of millions of users sending thousands of tokens each, and you arrive at the spend that built the data centers.

And remember the $O(n^2)$ from the math layer. Because attention compares every token to every other, a longer context costs compute that grows with the **square** of the token count — and the model must hold a **KV cache** whose memory grows *linearly* with every token in the conversation. That cache is a primary reason AI is starved for high-bandwidth memory. The token, in other words, is what runs straight into the **memory wall**.

So the chain is short and direct: **more tokens → more compute and more memory → more HBM and more data centers → the build-out**. The thing on this page is the atom that the whole [Circuit](#) is made of — see how it flows through the [memory supply chain](#) and the [economics in pictures](#).

The primary sources

Sennrich, Haddow & Birch (2016) — Neural Machine Translation of Rare Words with Subword Units · the paper that brought BPE to NLP.

Gage (1994) — A New Algorithm for Data Compression · the original byte-pair-encoding idea.

Radford et al. (2019) — GPT-2 · introduced byte-level BPE for language models.

OpenAI — tiktoken · the production tokenizer (cl100k_base, o200k_base).

Kudo & Richardson (2018) — SentencePiece · the language-agnostic alternative used by many open models.

Cite this chapter: Divergent Compute, "What is a token?", First Principles, 2026.
divergentcompute.com/first-principles-token · v1.0 · CC-BY.

The math the consensus is skipping — unlocked for life.

Join as a Founding Reader for free lifetime access to the interactive instruments, briefs, and data boards — before the gate closes.

500 LIFETIME FOUNDING SEATS · CLOSING AUG 25, 2026

Secure access

What is an embedding?

An embedding turns a token into a **point in space** — a list of numbers — arranged so that **meaning becomes distance**. Similar things sit close together, and relationships become directions you can do arithmetic with.

Read at your depth: [01 The answer](#) · [02 Intuition](#) · [03 Mechanics](#) · [04 The math](#) · [05 The code](#)
· [06 The economics](#) · [07 Sources](#)

Meaning, turned into geometry

In the [last chapter](#) a token became an integer ID. But an ID like 3404 carries no meaning — 3405 isn't "one more" of anything. So the very first thing a model does is replace each ID with an **embedding**: a vector of a few hundred to a few thousand numbers.

The magic is in *how* those numbers are arranged. The model learns them so that tokens used in similar ways end up **near each other** in the space. "Dog" and "cat" land close; "dog" and "Tuesday" land far apart. Meaning is no longer a symbol — it's a **location**.

And because it's a real vector space, you can do arithmetic on meaning. The famous example: take the vector for **king**, subtract **man**, add **woman** — and you land almost exactly on **queen**. The "royalty" and "gender" relationships are *directions* in the space. Try it:

THE EMBEDDING SPACE — EXPLORE IT

A 2-D schematic. Words that mean similar things cluster; relationships are parallel arrows.

● two
● one
● three

● wolf ● lion
● cat
● dog

● king

● queen

● prince

● princess

● tokyo

● japan

● man

● woman

● paris

● france

Show: king - man + woman

Highlight clusters

Real embeddings live in hundreds of dimensions; this is flattened to two so we can see it.

Where the numbers come from

Every model holds an **embedding table** — a big matrix with one row per vocabulary token. If the vocabulary has $|V|$ tokens and the model width is d , that table is $|V| \times d$ numbers (for GPT-2: ~50,000 rows of 768). Tokenizing gives you IDs; the embedding step is just a **lookup** — row 3404 of the table is the vector for that token.

Those numbers aren't hand-set; they're **learned**. During training the model nudges the vectors so that tokens appearing in similar contexts drift together — the old `word2vec` insight that "you shall know a word by the company it keeps." Modern LLMs learn their embeddings jointly with everything else, and they're *contextual*: the vector for "bank" shifts depending on whether the sentence is about rivers or money.

The same trick works on anything you can show a model — sentences, images, audio, code. Turn it into a vector, and "similar" becomes "close." That single idea is the engine under semantic search, recommendations, and retrieval.

Distance, similarity, and analogy

An embedding is a vector $\mathbf{v} \in \mathbb{R}^d$. To ask "how similar are two tokens," we don't use straight-line distance — we use the **angle** between their vectors, via **cosine similarity**:

$$\cos(\mathbf{u}, \mathbf{v}) = \frac{\mathbf{u} \cdot \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|} = \frac{\sum_{i=1}^d u_i v_i}{\sqrt{\sum_i u_i^2} \sqrt{\sum_i v_i^2}}$$

It runs from -1 (opposite) through 0 (unrelated) to 1 (identical direction). Angle, not magnitude, because what matters is *which way* a vector points in meaning-space, not how long it is.

Analogies fall out of the same geometry. If a relationship — "royal," say — is a consistent **direction** \mathbf{r} , then $\text{king} \approx \text{man} + \mathbf{r}$ and $\text{queen} \approx \text{woman} + \mathbf{r}$. Subtract to isolate the direction and re-add it elsewhere:

$$\text{king} - \text{man} + \text{woman} \approx \text{queen}$$

The model never "knew" that queens are royal women. The geometry encodes it, because that's how the words were used.

Similarity and analogy in a few lines

Cosine similarity and the king/queen analogy, on toy vectors — runnable as-is.

embeddings.py

```
import numpy as np

# toy 4-d embeddings (illustrative; real ones are hundreds of dims)
emb = {
    "king": np.array([0.9, 0.8, 0.1, 0.7]),
    "queen": np.array([0.9, 0.1, 0.8, 0.7]),
    "man": np.array([0.2, 0.8, 0.1, 0.1]),
    "woman": np.array([0.2, 0.1, 0.8, 0.1]),
}

def cosine(a, b):
    return a @ b / (np.linalg.norm(a) * np.linalg.norm(b))

# 1) similar words point the same way
print(round(cosine(emb["king"], emb["queen"]), 3)) # 0.749

# 2) analogy: king - man + woman ~ queen
v = emb["king"] - emb["man"] + emb["woman"]
best = max(emb, key=lambda w: cosine(v, emb[w]))
print(best, round(cosine(v, emb["queen"]), 3)) # queen 1.0
```

The interactive space above is the same idea, with the vectors arranged by hand so you can watch the parallelogram close.

Why "meaning as distance" became an industry

GEOMETRY → MONEY

The moment meaning is a location, finding relevant information becomes finding **nearby points**. That single move powers **semantic search**, recommendations, and **retrieval-augmented generation (RAG)** — the dominant way companies put their own data into an LLM. An entire category of infrastructure, the **vector database**, exists just to store billions of embeddings and find the nearest ones in milliseconds.

It isn't free. Every document, product, and message gets embedded (compute), and the vectors are held in fast memory to be searched at scale (memory). The embedding table itself is $|V| \times d$ parameters — the same $|V| \cdot d$ that the [token chapter](#) traded against sequence length. At the scale of the modern web, "turn everything into a vector and keep it searchable" is its own slice of the data-center demand.

So embeddings are the second atom — after the token — of the AI economy: the layer that turns *meaning* into something you can store, search, and bill for. See where it lands in the [Circuit](#).

The primary sources

Mikolov et al. (2013) — Efficient Estimation of Word Representations (word2vec) · the paper that made king – man + woman famous.

Pennington, Socher & Manning (2014) — GloVe · global word vectors from co-occurrence.

Devlin et al. (2018) — BERT · contextual embeddings: the same word, different vector by context.

Reimers & Gurevych (2019) — Sentence-BERT · embedding whole sentences for search.

Cite this chapter: Divergent Compute, "What is an embedding?", First Principles, 2026.
divergentcompute.com/first-principles-embedding · v1.0 · CC-BY.

The math the consensus is skipping — unlocked for life.

Join as a Founding Reader for free lifetime access to the interactive instruments, briefs, and data boards — before the gate closes.

500 LIFETIME FOUNDING SEATS · CLOSING AUG 25, 2026

Secure access

What is a neural network?

A neural network is a stack of **weighted sums**. Each unit multiplies its inputs by weights, adds them up, and passes the result through a simple bend. Learning is nothing more than **adjusting those weights** until the output is right.

Read at your depth: [01 The answer](#) · [02 Intuition](#) · [03 Mechanics](#) · [04 The math](#) · [05 The code](#)
· [06 The economics](#) · [07 Sources](#)

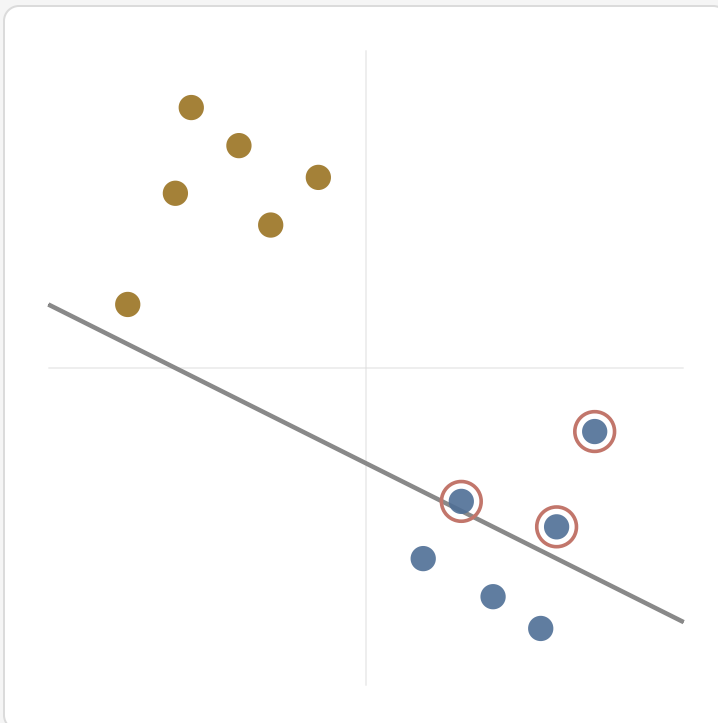
A neuron is just a line

Strip away the mystique and a single artificial neuron does one humble thing: it takes some numbers in, multiplies each by a **weight**, adds a **bias**, and decides yes-or-no. Geometrically, that "decide" is a **straight line** (a plane, in higher dimensions) splitting the input space in two.

Below is exactly one neuron with two inputs. Its weights set the angle and position of its decision line. **Drag the sliders** until the line cleanly separates the two groups of points — and watch the accuracy climb. That's a neuron "learning," done by hand.

ONE NEURON, THREE KNOBS

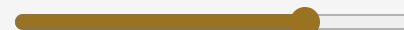
The line is where $w_1x + w_2y + b = 0$. Move the weights; separate the dots.



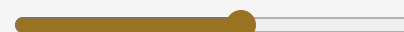
weight w_1 0.50



weight w_2 1.00



bias b 0.30



9/12 correct

One neuron draws one straight line. Stack a layer of them with a nonlinear bend between, and the lines compose into curves — that's how a deep network carves any shape.

From one neuron to a deep network

Three pieces turn that one line into something powerful:

- **Layers.** Put many neurons side by side (a layer), then feed their outputs into another layer, and another. Each layer's output is the next one's input — the **forward pass**.
- **The bend (activation).** Between layers sits a nonlinear function — `ReLU`, which just zeroes out negatives, is the workhorse. Without it, stacking layers would collapse back to a single line; *with* it, the network can bend its boundaries into arbitrary shapes. This is the whole reason depth helps.
- **Learning (gradient descent).** The network makes a prediction, measures how wrong it is (the **loss**), and computes which direction to nudge every weight to make the loss smaller — that direction is the **gradient**, found efficiently by **backpropagation**. Take a small step, repeat billions of times. The weights you dragged by hand, a network finds by rolling downhill.

That's the entire recipe. An LLM is this idea at staggering scale — the tokens and embeddings you've met flow through hundreds of such layers, and the "attention" of the next chapter is a particularly clever layer inside it.

Weighted sums, stacked

A single neuron with weights \mathbf{w} , bias b , and activation σ computes:

$$y = \sigma(\mathbf{w} \cdot \mathbf{x} + b) = \sigma(\sum_i w_i x_i + b)$$

A whole **layer** stacks many neurons, so the weights become a matrix W and the operation is a matrix-vector product:

$$\mathbf{h} = \sigma(W\mathbf{x} + \mathbf{b})$$

A deep network is just these **composed** — the output of one layer feeding the next:

$$f(\mathbf{x}) = \sigma(W_2 \sigma(W_1 \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2)$$

Training adjusts every weight by walking downhill on the loss L , scaled by a learning rate η :

$$W \leftarrow W - \eta \nabla_W L$$

That gradient $\nabla_W L$ is computed by backpropagation — the chain rule, applied layer by layer from the output back to the input. Everything else in modern AI is detail on top of this.

05 THE CODE

A forward pass in numpy

A two-layer network, front to back. Runnable as-is — this is the entire forward pass.

forward.py

```
import numpy as np

def relu(z):    return np.maximum(0, z)
def sigmoid(z): return 1 / (1 + np.exp(-z))

# layer 1: 3 inputs -> 4 hidden units ; layer 2: 4 -> 1 output
W1 = np.array([[ 0.5,-0.3, 0.8],
               [ 0.1, 0.9,-0.2],
               [-0.4, 0.2, 0.7],
               [ 0.3,-0.6, 0.1]])    # shape (4, 3)
b1 = np.array([0.0, 0.1, -0.2, 0.05])
W2 = np.array([[1.2, -0.7, 0.5, 0.9]])    # shape (1, 4)
b2 = np.array([0.05])

def forward(x):
    h = relu(W1 @ x + b1)    # hidden layer + the nonlinear bend
    y = sigmoid(W2 @ h + b2) # output as a probability
    return y

print(forward(np.array([1.0, 0.5, -1.0])))    # -> [0.36703]
```

The interactive above is the single-neuron version of the first line; a real model is millions of these matrices, learned rather than set.

Why "a stack of multiplies" costs a continent of power

MULTIPLIES → MONEY

A network's cost is set by one number: how many weighted multiplies it does. That's its **parameter count**, N . Running a model over one token takes roughly $2N$ **floating-point operations**; *training* it takes about $6ND$ — six times the parameters times the tokens it learns from. Those constants are not metaphors; they're the arithmetic that sizes the build-out.

So when a frontier model has $N \approx 10^{12}$ parameters and trains on $D \approx 10^{13}$ tokens, the training run is on the order of 10^{26} operations — months of a data center running flat out. The three sliders you moved are, at the frontier, a **trillion** of them, multiplied across trillions of tokens. That product *is* the demand for GPUs, power, and memory.

This is the hinge of the whole thesis: AI's compute bill is a direct function of network size, and size has only gone up. It's why "scaling laws" became an investment strategy — and why the [Circuit](#) exists. The humble weighted sum, multiplied enough times, built the data centers.

The primary sources

Rosenblatt (1958) — The Perceptron · the original artificial neuron.

Rumelhart, Hinton & Williams (1986) — Learning representations by back-propagating errors · backprop, the engine of learning.

LeCun, Bengio & Hinton (2015) — Deep Learning (Nature review) · the modern synthesis.

Kaplan et al. (2020) — Scaling Laws for Neural Language Models · where the 6ND compute rule comes from.

Cite this chapter: Divergent Compute, "What is a neural network?", First Principles, 2026. divergentcompute.com/first-principles-neural-network · v1.0 · CC-BY.

The math the consensus is skipping — unlocked for life.

Join as a Founding Reader for free lifetime access to the interactive instruments, briefs, and data boards — before the gate closes.

500 LIFETIME FOUNDING SEATS · CLOSES AUG 25, 2026

Secure access

What is attention?

Attention is the move that lets every token **look at every other token** and decide which ones matter for it right now. It's how a model figures out that "it" refers to "the animal" — and it's the heart of the transformer.

Read at your depth: [01 The answer](#) · [02 Intuition](#) · [03 Mechanics](#) · [04 The math](#) · [05 The code](#)
· [06 The economics](#) · [07 Sources](#)

Every word looks at every other word

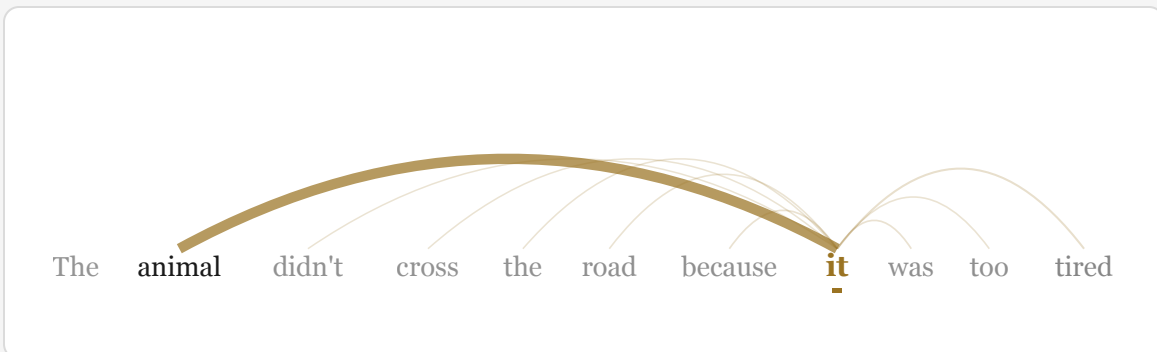
To understand a word, you need its context. "Bank" means something different by a river than in a sentence about money; "it" only makes sense once you know what it points to. Earlier models read left-to-right through a bottleneck and forgot.

Attention fixed that with a brutally direct idea: let every token look at *all* the others at once, and learn how much to weight each one.

Click a word below. The arcs show what it pays attention to — thicker means more. Notice where **"it"** looks:

ATTENTION — CLICK A WORD

A single attention head over one sentence. Arcs show how much each word attends to the others.



"it" attends most to **"animal"** (67%). That's how the model resolves what "it" refers to.

Illustrative weights for one head. Real models run dozens of these heads in parallel, each learning a different kind of relationship.

Query, key, value

Each token produces three vectors from its embedding, via learned weight matrices:

- **Query** — what this token is looking for.
- **Key** — what each token offers, as an advertisement.
- **Value** — the actual content each token will hand over if attended to.

The mechanism is a soft dictionary lookup. Compare one token's **query** against every token's **key** (a dot product) to get a relevance score; the better the match, the higher the score. Run those scores through a `softmax` so they become weights that sum to one, then take a **weighted average of the values**. That blend is the token's new, context-aware representation — "it" has now mixed in a large dose of "animal."

Stack this inside the network layers, run many heads in parallel (each catching a different relationship — syntax, coreference, topic), and repeat over dozens of layers. That stack is the **transformer**, and it's the next chapter.

Scaled dot-product attention

Pack the queries, keys and values for all n tokens into matrices Q, K, V , each of shape $n \times d_k$. The entire operation is one line:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^\top}{\sqrt{d_k}}\right) V$$

Read it inside-out. QK^\top is an $n \times n$ matrix — **every token's query dotted with every token's key**, the full grid of relevance scores. Dividing by $\sqrt{d_k}$ keeps the numbers from blowing up as dimension grows. The `softmax` turns each row into a probability distribution (the attention weights). Multiplying by V takes the weighted average of the values.

Keep your eye on that QK^\top . It is $n \times n$ — its size grows with the **square** of the number of tokens. Hold that thought; it is the most expensive fact in all of AI, and it returns in layer 06.

Attention in seven lines

The whole mechanism, runnable. Note the attention weights — each row sums to one.

attention.py

```
import numpy as np

def softmax(x):
    e = np.exp(x - x.max(axis=-1, keepdims=True))
    return e / e.sum(axis=-1, keepdims=True)

def attention(Q, K, V):
    d_k = Q.shape[-1]
    scores = Q @ K.T / np.sqrt(d_k) # (n, n): every query vs every key
    weights = softmax(scores) # each row sums to 1
    return weights @ V, weights # blended values, and the attention map

np.random.seed(0)
Q = np.random.randn(3, 4); K = np.random.randn(3, 4); V = np.random.randn(3, 4)
out, w = attention(Q, K, V)
print(w.round(2)) # [[0.68 0.3 0.02]
                  # [0.29 0.7 0.01]
                  # [0.5 0.19 0.31]]
print(w.sum(axis=1)) # [1. 1. 1.]
```

The square that built the data centers

$N^2 \rightarrow$ MONEY

Attention's superpower — letting every token see every other — is also its bill. That QK^T grid is $n \times n$, so the compute grows as $O(n^2)$ in the number of tokens.

Double the context, and you quadruple the attention work. This single fact is the most consequential number in AI economics.

It's why early models capped context at a few thousand tokens, why "1 million token context" is a genuine engineering feat rather than a setting, and why running models is **memory-bound**: to avoid recomputing, models cache every token's keys and values — the **KV cache** — whose size grows with the sequence and devours high-bandwidth memory. The n you met in the [token chapter](#), *attention squares*.

So the line of demand runs straight through this layer: longer, smarter context \rightarrow quadratic compute and linear-but-relentless memory \rightarrow more GPUs and more **HBM** \rightarrow the build-out. A whole research frontier (FlashAttention, sparse and linear attention) exists just to soften that exponent. When the [Circuit](#) talks about the memory wall, this n^2 is the wall.

The primary sources

Vaswani et al. (2017) — Attention Is All You Need · the transformer, and this exact formula.

Bahdanau, Cho & Bengio (2014) — Neural Machine Translation by Jointly Learning to Align and Translate · attention, before transformers.

Dao et al. (2022) — FlashAttention · making the $O(n^2)$ memory-efficient on real GPUs.

Alammar — The Illustrated Transformer · the canonical visual walkthrough.

Cite this chapter: Divergent Compute, "What is attention?", First Principles, 2026.
divergentcompute.com/first-principles-attention · v1.0 · CC-BY.

The math the consensus is skipping — unlocked for life.

Join as a Founding Reader for free lifetime access to the interactive instruments, briefs, and data boards — before the gate closes.

500 LIFETIME FOUNDING SEATS · CLOSING AUG 25, 2026

Secure access

What is a transformer?

A transformer is the **assembly**. Take embeddings, attention, and a small neural network, wire them into a repeatable **block**, and stack that block dozens of times. That stack is what an LLM is.

Read at your depth: [01 The answer](#) · [02 Intuition](#) · [03 Mechanics](#) · [04 The math](#) · [05 The code](#)
· [06 The economics](#) · [07 Sources](#)

The parts you already know, wired together

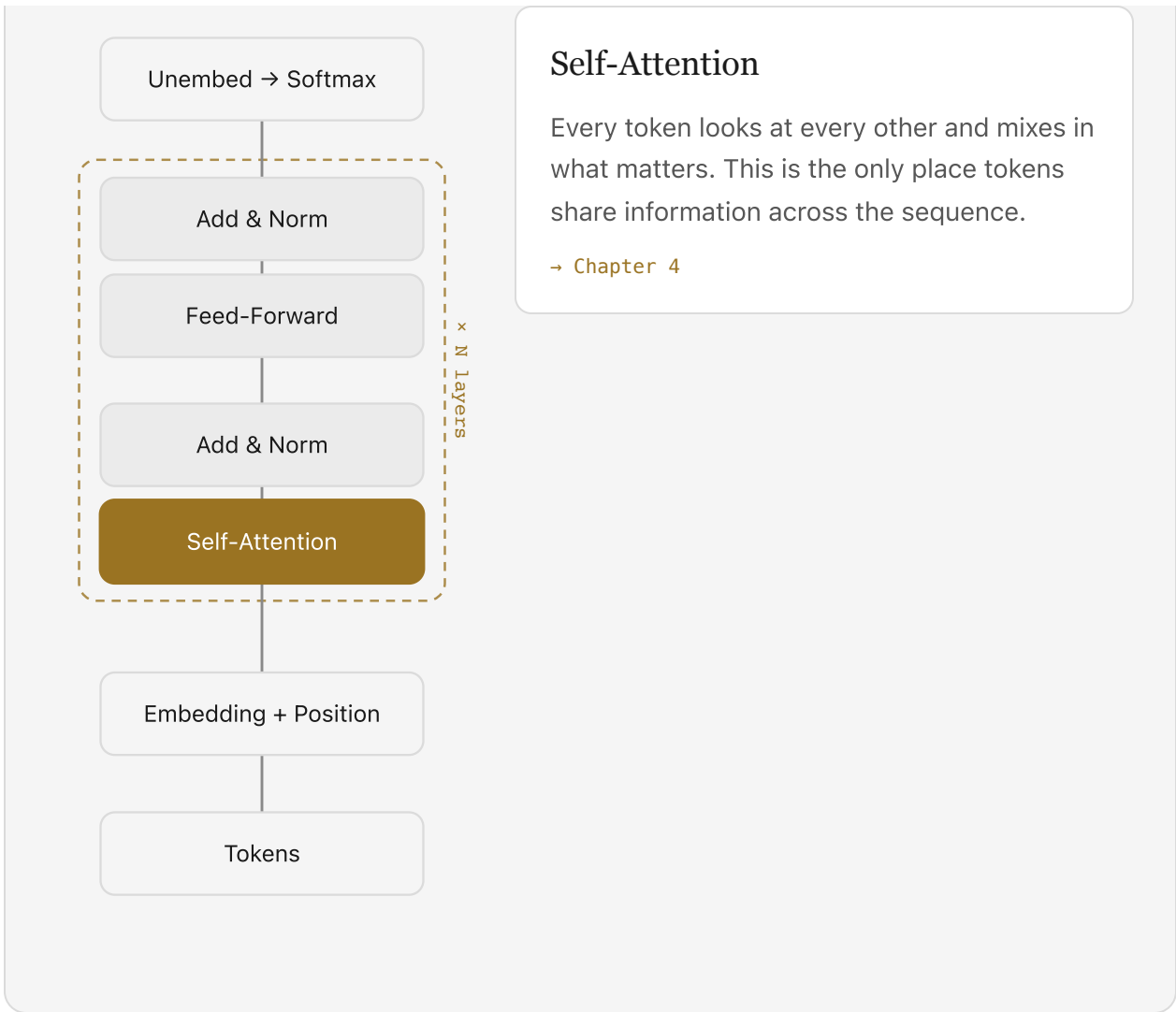
You've met every piece. A transformer just connects them. Text becomes tokens, tokens become embeddings, and then the data flows through a stack of identical **blocks**. Inside each block, attention lets tokens share information, and a small neural network processes each one. Repeat, and out the top come probabilities for the next token.

Click through the architecture — each stage is one of the earlier chapters doing its job:

THE TRANSFORMER — CLICK ANY STAGE

Data flows up. The shaded block is the unit that repeats.





Self-Attention

Every token looks at every other and mixes in what matters. This is the only place tokens share information across the sequence.

→ Chapter 4

One block, repeated

The whole architecture is one block applied over and over. Inside it, two ideas you've seen, plus two pieces of plumbing that make deep stacks trainable:

- **Self-attention** mixes information *across* tokens — the only place they talk to each other.
- **The feed-forward network** then processes each token *independently*. This is where most of a model's parameters actually live.
- **Residual connections** add a layer's input back to its output, so information (and gradients) can skip straight through a hundred layers without vanishing.
- **Layer normalization** rescales the numbers at each step to keep training stable.

One detail matters: attention is **order-blind** — shuffle the tokens and it gives the same answer. So before the first block, each embedding gets a **positional encoding** added, a signal telling the model where each token sits. Stack ~12 blocks for a small model, ~100 for a frontier one, cap it with an output layer that turns the final vectors into next-token probabilities, and you have GPT.

The block as two residual steps

Let \mathbf{x} be the matrix of token vectors entering a block. The block is exactly two residual-and-normalize steps — first attention, then the feed-forward network:

$$\mathbf{x}' = \text{LayerNorm}(\mathbf{x} + \text{Attention}(\mathbf{x}))$$

$$\mathbf{x}'' = \text{LayerNorm}(\mathbf{x}' + \text{FFN}(\mathbf{x}'))$$

where the feed-forward network is a two-layer MLP applied to each token, $\text{FFN}(\mathbf{x}) = \max(0, \mathbf{x}W_1)W_2$. Crucially the block maps $\mathbb{R}^{n \times d} \rightarrow \mathbb{R}^{n \times d}$ — **same shape in, same shape out** — which is precisely why you can stack it N times:

$$f(\mathbf{x}) = \text{Block}_N(\cdots \text{Block}_2(\text{Block}_1(\mathbf{x})))$$

The final vectors are projected by an unembedding matrix to vocabulary-sized scores (logits) and softmaxed into the probability of the next token. That's the entire forward pass of a GPT.

A transformer block in numpy

The block, composing the pieces from the earlier chapters. Note the output shape equals the input shape — that's the stacking property.

block.py

```
import numpy as np
np.random.seed(0)

def softmax(x):
    e = np.exp(x - x.max(axis=-1, keepdims=True)); return e / e.sum(axis=-1, keepdims=True)
def attention(x, Wq, Wk, Wv):
    Q, K, V = x@Wq, x@Wk, x@Wv
    return softmax(Q@K.T / np.sqrt(Q.shape[-1])) @ V
def layernorm(x):
    return (x - x.mean(-1, keepdims=True)) / (x.std(-1, keepdims=True) + 1e-5)
def ffn(x, W1, W2):
    return np.maximum(0, x@W1) @ W2 # 2-layer ReLU MLP

def block(x, Wq, Wk, Wv, W1, W2):
    x = layernorm(x + attention(x, Wq, Wk, Wv)) # attention + residual
    x = layernorm(x + ffn(x, W1, W2)) # feed-forward + residual
    return x

x = np.random.randn(3, 4) # 3 tokens, width 4
Wq, Wk, Wv = (np.random.randn(4, 4) for _ in range(3))
W1, W2 = np.random.randn(4, 8), np.random.randn(8, 4) # FFN hidden = 8
print(block(x, Wq, Wk, Wv, W1, W2).shape) # (3, 4) - stack it
```

Depth times width times the square

ARCHITECTURE → MONEY

The transformer's bill is its shape. Each of the N blocks runs an attention (the $O(n^2)$ from the [last chapter](#)) and a feed-forward network whose cost scales with the square of the width, $O(n d^2)$. Multiply by the number of layers and you have the whole forward pass. A frontier model is roughly $N \approx 100$ blocks, width d in the tens of thousands, over n in the thousands — and it runs that for every token, for hundreds of millions of users.

This is why the architecture *is* the capex. Every knob that makes models better — more layers, more width, longer context — multiplies straight into FLOPs, GPUs, power, and memory. The transformer turned "make it bigger" into a reliable recipe, and that recipe is what the entire build-out is paying for.

You now hold the whole chain: a [token](#) (n), its [embedding](#) (d), the [attention](#) (n^2), the [network](#) (N layers of weights) — assembled here into the machine whose cost is the [Circuit](#). The transformer is where the foundations become an economy.

The primary sources

Vaswani et al. (2017) — Attention Is All You Need · the transformer architecture itself.

Radford et al. (2018) — GPT · the decoder-only transformer that became the LLM.

Phuong & Hutter (2022) — Formal Algorithms for Transformers · the architecture written out precisely.

Alammar — The Illustrated Transformer · the canonical visual walkthrough.

Cite this chapter: Divergent Compute, "What is a transformer?", First Principles, 2026.
divergentcompute.com/first-principles-transformer · v1.0 · CC-BY.

The math the consensus is skipping — unlocked for life.

Join as a Founding Reader for free lifetime access to the interactive instruments, briefs, and data boards — before the gate closes.

500 LIFETIME FOUNDING SEATS · CLOSING AUG 25, 2026

Secure access

Parameters & weights

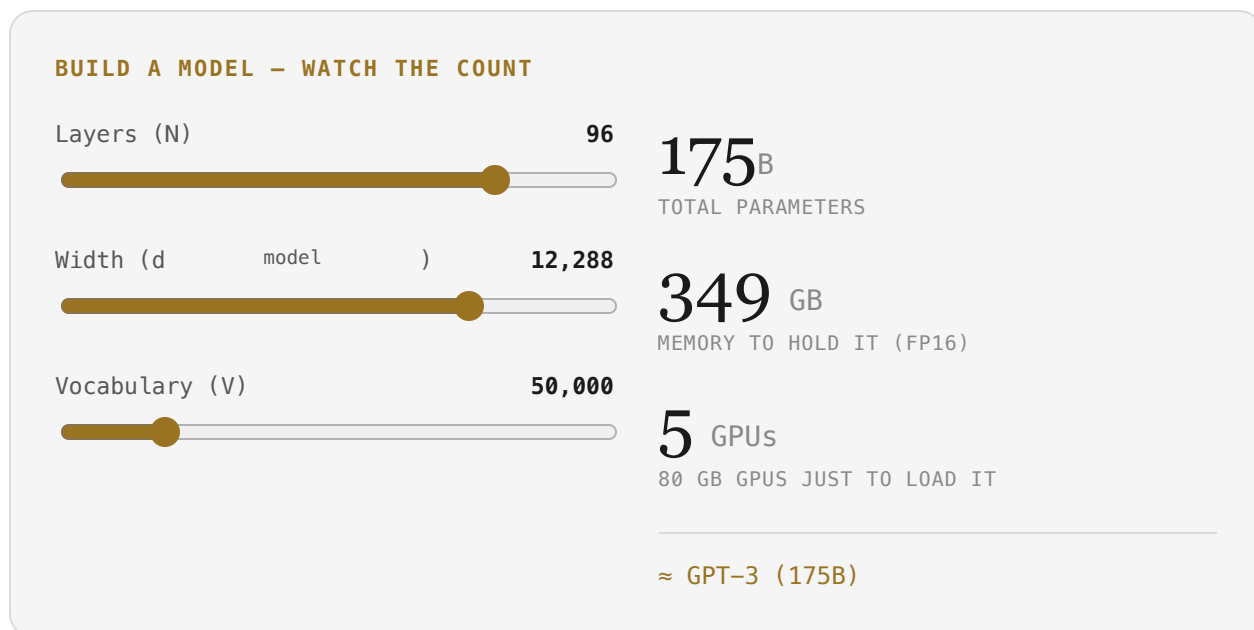
A parameter is a single **learned number** — one weight in one of the model's matrices. A large model is billions of them, and together those numbers **are** the model. "175 billion parameters" means 175 billion learned knobs.

Read at your depth: [01 The answer](#) · [02 Intuition](#) · [03 Mechanics](#) · [04 The math](#) · [05 The code](#)
· [06 The economics](#) · [07 Sources](#)

The model is its numbers

Every multiply you've met — the embedding lookup, the attention projections, the feed-forward layers — uses weights. A **parameter** is just one of those weights: a single number the model learned during training. There's nothing else in the model. The "knowledge" isn't stored anywhere special; it's distributed across all the weights, and the model file you download is literally a giant list of them.

That's why parameter count is the headline spec. It sets how much the model can learn, how much memory it takes to hold, and how much compute it takes to run. Build your own and watch those three numbers move — set the depth and width of a transformer:



Where the billions live

The parameters aren't spread evenly. In each transformer layer they sit in a handful of matrices:

- **Attention** — four square matrices (W_Q, W_K, W_V, W_O), each $d \times d$, so $4d^2$ per layer.
- **Feed-forward** — two matrices that expand to a wider hidden size (usually $4d$) and back, about $8d^2$ per layer. **This is where most parameters live.**
- **Embeddings** — one $V \times d$ table at the bottom, and a same-sized one at the top (often tied to save space).
- **Norms & biases** — a rounding error by comparison.

So a layer holds roughly $12d^2$ parameters, and a model with N layers holds about $12Nd^2$ plus the embeddings. They begin as random noise and are nudged, one tiny gradient step at a time, until the whole pile of numbers predicts text well. Training is the act of *setting* these parameters; everything else in AI is what you do with them once they're set.

Counting the knobs

With N layers, width d , and vocabulary V , the parameter count is dominated by the layers:

$$N_{\text{params}} \approx \underbrace{N(4d^2 + 8d^2)}_{\text{attention + FFN}} + \underbrace{Vd}_{\text{embedding}} = 12Nd^2 + Vd$$

For a model of GPT-3's shape — $N = 96$, $d = 12288$, $V \approx 50,257$ — that comes to about 175 billion, which is exactly the famous number. Notice the d^2 : **widening a model is quadratically expensive in parameters**, while adding layers is only linear.

Two consequences set the cost. The memory to hold the weights is N_{params} times the bytes per number — 2 bytes in half-precision, 1 in int8. And running the model takes about $2N_{\text{params}}$ floating-point operations *per token*. So the single number on the box determines both the storage and the compute.

The parameter calculator, in code

The exact function behind the slider above — runnable, and it reproduces GPT-3's 175B.

params.py

```
def transformer_params(n_layers, d_model, vocab, d_ff=None):
    d_ff = d_ff or 4 * d_model
    per_layer = 4 * d_model**2 + 2 * d_model * d_ff    # attention (QKV0) + FFN
    embedding = vocab * d_model                        # token embedding table
    return n_layers * per_layer + embedding

# GPT-3 scale: 96 layers, width 12288, vocab 50257
N = transformer_params(96, 12288, 50257)
print(f"{N:,} parameters")                          # 174,563,733,504 parameters
print(f"{N/1e9:.0f}B params | {N*2/1e9:.0f} GB in fp16") # 175B params | 340GB
```

The number on the box is the bill

COUNT → MONEY

Parameter count is the closest thing the AI economy has to a single price tag. It sets **memory** — a 175B model needs about **350 GB** just to sit in half-precision, about five 80 GB GPUs before it has done any work. It sets **compute** — roughly $2N$ operations per token to run, and about $6N \cdot D$ to train. And by the **scaling laws**, it sets **capability**: bigger has reliably meant better, which is why the number keeps climbing.

That is the whole flywheel of the build-out. "Make it bigger" improves the product, but every extra parameter is more HBM to hold it, more GPUs to serve it, and more energy to train it. The race to larger N is the race that fills the data centers — and the reason memory, not raw compute, is so often the binding constraint.

So when you read "405 billion parameters," translate it: that many learned numbers, that much memory, that much silicon. The parameter count you just dialed is the unit in which the entire [Circuit](#) is denominated.

The primary sources

Brown et al. (2020) — Language Models are Few-Shot Learners (GPT-3) · the 175-billion-parameter model.

Kaplan et al. (2020) — Scaling Laws for Neural Language Models · why bigger reliably means better.

Hoffmann et al. (2022) — Training Compute-Optimal LLMs (Chinchilla) · how to balance parameters against training tokens.

Puong & Hutter (2022) — Formal Algorithms for Transformers · the parameter tables, written out.

Cite this chapter: Divergent Compute, "Parameters & weights", First Principles, 2026.
divergentcompute.com/first-principles-parameters · v1.0 · CC-BY.

The math the consensus is skipping — unlocked for life.

Join as a Founding Reader for free lifetime access to the interactive instruments, briefs, and data boards — before the gate closes.

500 LIFETIME FOUNDING SEATS · CLOSING AUG 25, 2026

Secure access

The context window

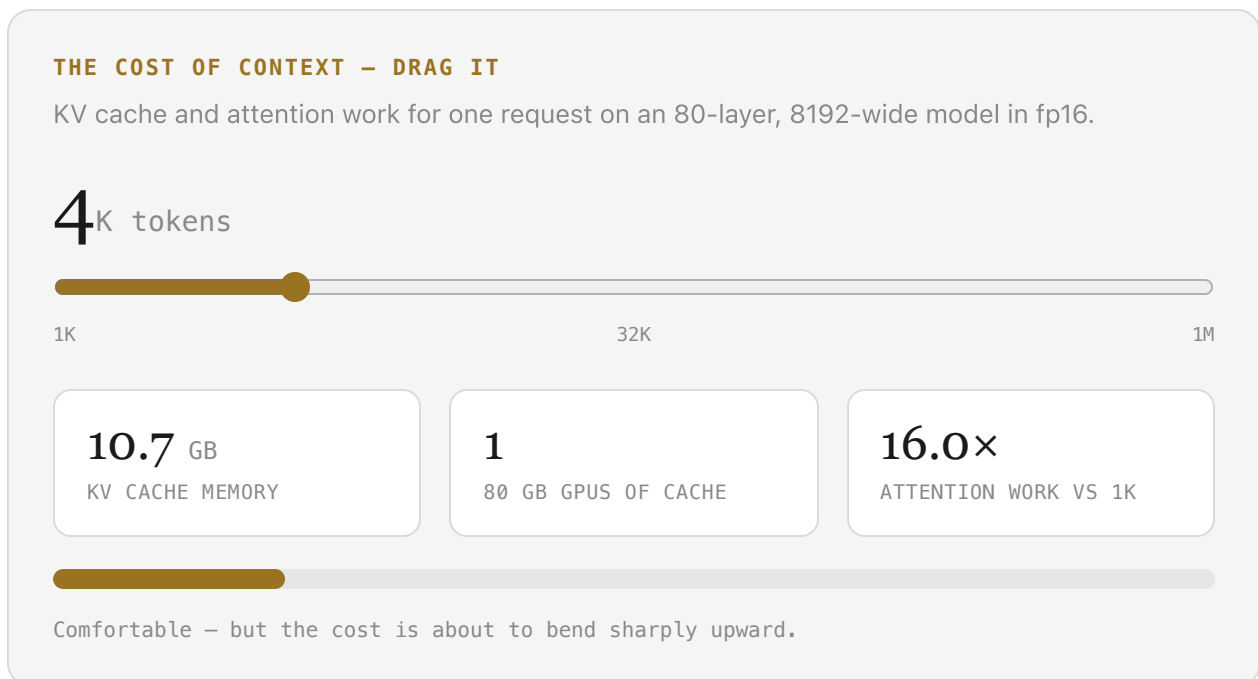
The context window is a model's **working memory** — the number of tokens it can hold and attend to at once. Everything inside it the model can use; everything outside it, it never sees. And making it bigger is brutally expensive.

Read at your depth: [01 The answer](#) · [02 Intuition](#) · [03 Mechanics](#) · [04 The math](#) · [05 The code](#)
· [06 The economics](#) · [07 Sources](#)

How much the model can hold at once

A model doesn't have long-term memory between turns; within a single request it has the **context window** — the maximum number of tokens it can read and reason over together. A whole conversation, a pasted document, the system instructions: it all has to fit. Run past the limit and the oldest tokens fall off the edge, or get truncated, and the model simply doesn't know they existed.

Windows have grown fast — from ~2,000 tokens in early GPT-3 to 128,000 in GPT-4-class models to over a million in the largest. But that growth fights a hard wall you've already met: attention's cost grows with the **square** of the window, and the model must cache every token's keys and values in fast memory. Drag the window and watch what it costs (on a 70B-class model):



Why it's capped, and how it's stretched

Two costs set the ceiling, both from attention:

- **Compute grows quadratically.** Every token attends to every other, so the work scales as n^2 . Going from a 1K window to a 128K one is a 128x increase in length — and therefore about **16,000x** the attention work (128 squared). It explodes.
- **Memory grows linearly — and ruinously.** To avoid recomputing, the model stores every token's keys and values, the **KV cache**. That cache grows with the window *and* the model's depth and width, and it lives in scarce high-bandwidth memory.

So engineers fight on three fronts: **position tricks** like RoPE and ALiBi that let a model generalize to lengths it wasn't trained on; **efficient attention** like FlashAttention, sliding-window, and sparse patterns that cut the constants (or even the exponent); and simply **more memory**. One caveat worth knowing: even when a model *can* read a million tokens, it often attends less to the middle — the "lost in the middle" effect — so a bigger window isn't automatically better understanding.

Quadratic compute, linear cache

For a window of n tokens, width d , and N layers, the two costs are:

$$\text{attention compute} \propto n^2 d \qquad \text{KV cache bytes} = 2 N n d b$$

The 2 counts keys and values, b is the bytes per number (2 in fp16). The compute term is the famous quadratic: doubling n quadruples the attention work. The cache term is "only" linear in n — but the constant is enormous, because it also multiplies by every layer and the full width.

Make it concrete. For an 80-layer, 8192-wide model, each token costs $2 \cdot 80 \cdot 8192 \cdot 2 \approx 2.6$ MB of cache. So a 128K-token window holds about **340 GB** of KV cache — roughly five 80 GB GPUs, before the model produces a single word. A one-million-token window is measured in **terabytes**. (These are full multi-head-attention figures — the honest worst case. Real 70B-class models use **grouped-query attention**, sharing keys and values across heads to shrink the cache by roughly 8x; it is precisely the trick invented to fight this wall.) That is why long context is a hardware problem, not a software setting.

The KV-cache calculator

The exact function behind the slider — runnable.

kv_cache.py

```
def kv_cache_gb(n_tokens, n_layers, d_model, bytes_per=2):
    # 2 = one key + one value vector per token, per layer
    return 2 * n_layers * n_tokens * d_model * bytes_per / 1e9

# an 80-layer, 8192-wide model in fp16
for ctx in [4096, 32768, 131072, 1_000_000]:
    print(f"{ctx:>9,} tokens -> {kv_cache_gb(ctx, 80, 8192):8.1f} GB of KV cache")
# ->    4,096 tokens ->    10.7 GB of KV cache
#      32,768 tokens ->    85.9 GB of KV cache
#     131,072 tokens ->   343.6 GB of KV cache
#    1,000,000 tokens -> 2621.4 GB of KV cache
```

Where the square meets the customer

WINDOW → MONEY

The context window is where every cost in this Part comes due at once. The [tokens](#) set n ; [attention](#) squares it for compute; the KV cache turns it into **hundreds of gigabytes of high-bandwidth memory** per long request. This is the most direct reason inference is memory-bound, and why providers charge more for long-context calls — you are renting scarce HBM by the token-squared.

It also reframes the build-out. A huge share of data-center spending isn't about training ever-larger models; it's about **servicing** them — holding millions of users' KV caches in memory at once. Every leap in advertised context length ripples straight into demand for the exact resource the [Circuit](#) calls the memory wall.

So the context window is the perfect closing note for the foundations: it is the single setting where the token, the embedding, attention, the network, the transformer, and the parameter count all collapse into one number you pay for. Read the spec "1M token context," and now you know precisely what it costs to keep that promise.

That completes Part I — the Foundations. From a [token](#) to the full [transformer](#) and the economics of running it, you now have the whole machine in view. Part II opens the model itself: how LLMs are trained, fine-tuned, and how they differ. [Back to the curriculum →](#)

The primary sources

Su et al. (2021) — RoFormer (RoPE) · rotary position encodings that help models extend their length.

Dao et al. (2022) — FlashAttention · making long-context attention fit in memory.

Liu et al. (2023) — Lost in the Middle · why a longer window isn't always better understanding.

Beltagy et al. (2020) — Longformer · sparse attention for long documents.

Cite this chapter: Divergent Compute, "The context window", First Principles, 2026.
divergentcompute.com/first-principles-context-window · v1.0 · CC-BY.

The math the consensus is skipping — unlocked for life.

Join as a Founding Reader for free lifetime access to the interactive instruments, briefs, and data boards — before the gate closes.

500 LIFETIME FOUNDING SEATS · CLOSES AUG 25, 2026

Institutional email

Secure access

What is an LLM?

A large language model is a **transformer** trained to do one thing: **predict the next token**. Run that prediction over and over, feeding each guess back in, and the next-token machine becomes a writer.

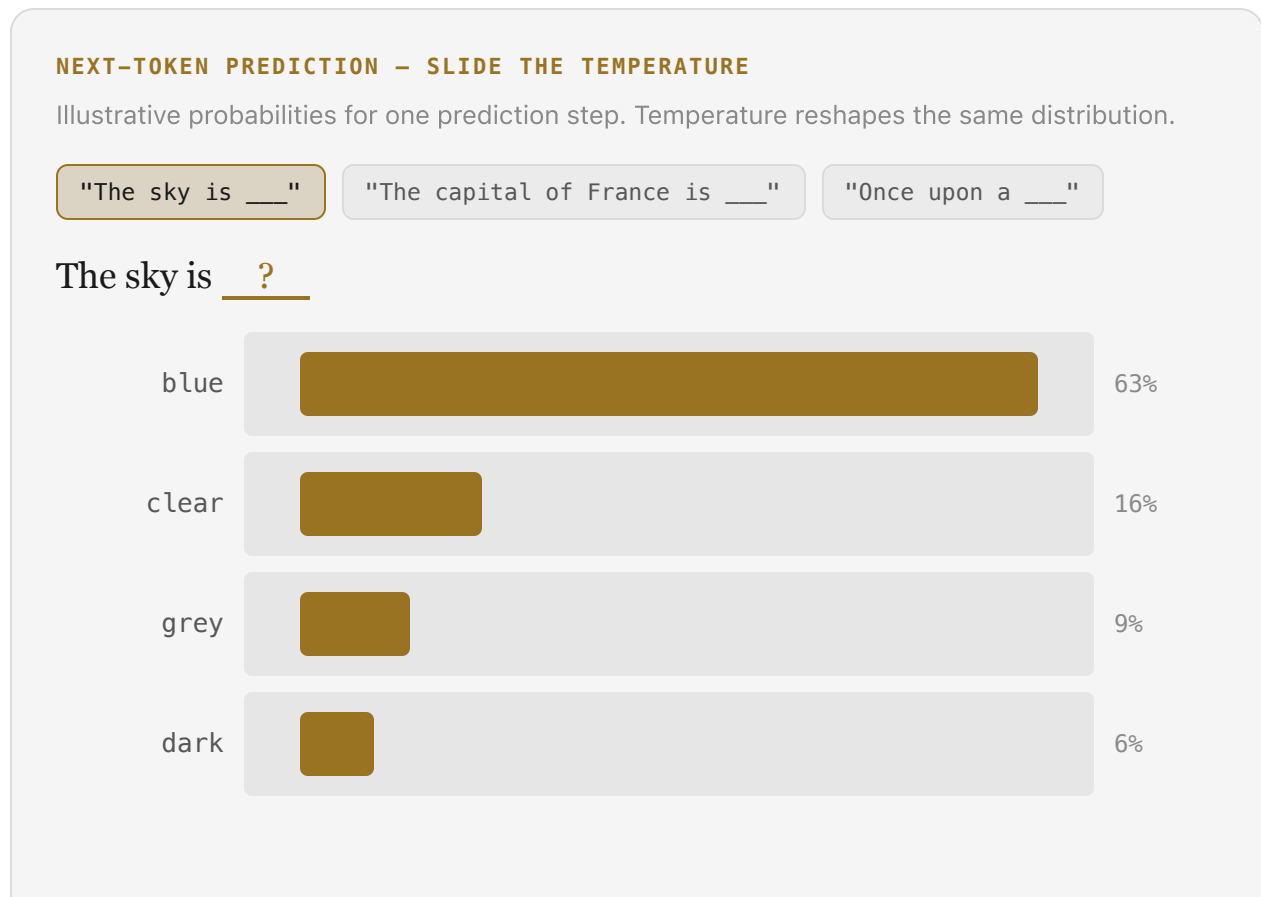
Read at your depth: [01 The answer](#) · [02 Intuition](#) · [03 Mechanics](#) · [04 The math](#) · [05 The code](#)
· [06 The economics](#) · [07 Sources](#)

It is autocomplete — taken extremely seriously

Everything in Part I builds to this. An LLM takes your text as tokens, runs it through the transformer stack, and produces a probability for *every possible next token*. It picks one, sticks it on the end, and runs again on the longer text. Repeat until done. That loop — **predict, append, repeat** — is all "generation" is.

The surprise of the last few years is that doing this one humble task well enough, at enough scale, produces translation, code, reasoning, and conversation as side effects. To answer a question, the most likely continuation of "the answer is" turns out to be the answer.

Below is that prediction step. The model gives a distribution over next tokens; a single knob, the **temperature**, controls how boldly it picks. Slide it:



the



3%

falling



2%

temperature



1.00

Temperature 1.0 is the model's raw distribution. Drag it down for safe, up for creative.

Predict, sample, repeat

The mechanics are a short loop:

- **The objective.** During training the model is shown oceans of text and asked, at every position, to predict the next token. Its only goal is to make the real next token likely. No labels, no human in the loop — just "guess what comes next," billions of times.
- **The output.** At generation time the final vector becomes a **logit** for every token in the vocabulary; a `softmax` turns those into probabilities.
- **Sampling.** You then pick a token. **Greedy** takes the most likely. **Temperature** scales the logits to make the choice sharper or flatter. **Top-p** samples only from the smallest set of tokens that covers, say, 90% of the probability. These knobs are the difference between a dull, deterministic model and a lively one.
- **Autoregression.** Append the chosen token and run the whole thing again. Each new token is conditioned on everything so far — which is exactly why the context window matters.

So an "LLM" is not a database of answers. It is a single learned function for "what token is likely next," wrapped in a loop. The next chapters open up where that function's knowledge comes from — pretraining — and how it's shaped to be helpful.

A distribution over the vocabulary

Given the tokens so far, the model outputs a logit z_i for each token i in the vocabulary. With temperature T , the probability of the next token is:

$$P(\text{token} = i) = \frac{\exp(z_i/T)}{\sum_j \exp(z_j/T)}$$

As $T \rightarrow 0$ this collapses onto the single highest-logit token (greedy, deterministic); at $T = 1$ it is the model's raw distribution; as T grows the probabilities flatten toward uniform — more random, eventually nonsense. Generating a whole sequence just multiplies these step by step:

$$P(t_1, t_2, \dots, t_m) = \prod_{k=1}^m P(t_k | t_1, \dots, t_{k-1})$$

Training maximizes exactly this probability on real text — equivalently, it minimizes the average **cross-entropy** (the negative log of the right token's probability).

"Predict the next token" is the entire objective.

Sampling, and the generation loop

Temperature sampling, then the autoregressive loop. Runnable — it's the whole idea.

generate.py

```
import numpy as np

def softmax(z, T=1.0):
    z = np.array(z, dtype=float) / T
    e = np.exp(z - z.max()); return e / e.sum()

# illustrative logits over a tiny vocabulary, after "The sky is"
vocab = ["blue", "clear", "grey", "dark", "falling", "the"]
logits = [4.2, 2.8, 2.3, 1.9, 0.8, 1.2]

for T in [0.5, 1.0, 1.7]:
    p = softmax(logits, T)
    print(f"T={T}:", ", ".join(f"{w} {pi:.2f}" for w, pi in zip(vocab, p)))
# T=0.5: blue 0.91, clear 0.06, grey 0.02, dark 0.01, falling 0.00, the 0.00
# T=1.0: blue 0.63, clear 0.16, grey 0.09, dark 0.06, falling 0.02, the 0.03
# T=1.7: blue 0.43, clear 0.19, grey 0.14, dark 0.11, falling 0.06, the 0.07

def generate(next_logits, prompt, n, T=1.0):    # next_logits: context -> log
    toks = list(prompt)
    for _ in range(n):
        p = softmax(next_logits(toks), T)
        toks.append(int(np.random.choice(len(p), p=p)))    # sample the next to
    return toks
```

Sold by the token, one at a time

PREDICTION → MONEY

Because generation is a loop, an LLM cannot produce its answer all at once — it must run a **full forward pass for every single token** it writes, each one conditioned on all the tokens before it. That sequential dependency is why responses stream in word by word, and why latency and cost scale with output length. Every token is one trip through $\sim 2N$ [parameters](#).

This is the unit the whole industry is priced in. Providers bill per input and output token; the output tokens cost more precisely because each one is a fresh, unparallelizable forward pass. Multiply by hundreds of millions of users generating billions of tokens a day, and "predict the next token" becomes the single largest recurring workload in computing — the **inference** demand that, even more than training, is what the data centers are being built to serve.

So the humble loop you just sped through is the business. The token is the product; the model is the factory; and the [Circuit](#) is the account of whether that factory's output ever justifies its cost.

The primary sources

Radford et al. (2019) — GPT-2 · language modeling as a path to general capability.

Brown et al. (2020) — GPT-3 · scale turns next-token prediction into few-shot learning.

Holtzman et al. (2019) — The Curious Case of Neural Text Degeneration · why top-p / nucleus sampling beats greedy.

Jurafsky & Martin — Speech and Language Processing · the language-modeling objective, from the ground up.

Cite this chapter: Divergent Compute, "What is an LLM?", First Principles, 2026.
divergentcompute.com/first-principles-llm · v1.0 · CC-BY.

The math the consensus is skipping — unlocked for life.

Join as a Founding Reader for free lifetime access to the interactive instruments, briefs, and data boards — before the gate closes.

500 LIFETIME FOUNDING SEATS · CLOSING AUG 25, 2026

Institutional email

Secure access

How models are pretrained

Pretraining is the long, brutally expensive first phase where a model reads **trillions of tokens** of text and, by predicting the next one over and over, teaches itself language, facts, and reasoning — with no labels and no human in the loop.

Read at your depth: [01 The answer](#) · [02 Intuition](#) · [03 Mechanics](#) · [04 The math](#) · [05 The code](#)
· [06 The economics](#) · [07 Sources](#)

Learning the world by guessing the next word

The model from the [last chapter](#) isn't born knowing anything — its [parameters](#) start as random noise. Pretraining is how those billions of numbers become a model of language. The recipe is almost absurdly simple: feed it a colossal pile of text, and at every position ask it to predict the next token. When it's wrong, nudge the weights to be a little less wrong. Do that for trillions of tokens.

Because the right answer is always just "the next word in the text," no human has to label anything — the data **supervises itself**. To predict text well, the model is quietly forced to learn grammar, facts, code, arithmetic, and reasoning, because all of those help it guess what comes next. Drag through a training run and watch capability emerge:

A TRAINING RUN — DRAG THROUGH IT

Illustrative. The loss falls and the model's writing sharpens as it sees more tokens.

0.4T

TRAINING TOKENS SEEN

7.85

LOSS (LOWER IS BETTER)



OUTPUT · RANDOM WEIGHTS

“qx z* the the aa . the the the the”

02 MECHANICS

The data, the loss, and the months of compute

- **The data.** A frontier model trains on something like 10–15 *trillion* tokens — a filtered, deduplicated slice of the web, books, and code. Data quality and mix matter as much as quantity; much of the craft of pretraining is in the cleaning.
- **The objective.** At each position the model predicts a distribution over the next token, and the loss is the **cross-entropy** — how surprised it was by the real next token. Averaged over the whole corpus, that single number is what training drives down.
- **The compute.** Each token is processed by every parameter; training cost is about $6ND$ — six times the parameters times the tokens. For a frontier model that is on the order of 10^{25} – 10^{26} operations: **thousands of GPUs running for months.**
- **Emergence.** As scale grows the loss falls along a smooth power law, but specific abilities — arithmetic, translation, in-context learning — can appear rather suddenly once the model is big enough. You can't fully predict *what* a bigger model will be able to do, only that it will do *better*.

The output is a **base model**: fluent, knowledgeable, and completely unhelpful. It will happily continue your text, but it hasn't learned to follow instructions or be safe.

Turning it into a usable assistant is the next chapter.

Cross-entropy, minimized over a corpus

For one position, with the model's predicted distribution P and the true next token t , the loss is the negative log-probability it assigned to the truth:

$$\ell = -\log P(t)$$

Pretraining minimizes the average of this over the whole corpus of D tokens — the **cross-entropy loss**:

$\mathcal{L} = -\frac{1}{D} \sum_{k=1}^D \log P(t_k \mid t_{< k})$
 by gradient descent on the parameters θ : $\theta \leftarrow \theta - \eta \nabla_{\theta} \mathcal{L}$. Empirically the loss follows a **scaling law** — it falls as a power of the compute C :

$$\mathcal{L}(C) \approx \mathcal{L}_{\infty} + \left(\frac{C_0}{C}\right)^{\alpha}$$

This smooth, predictable curve is what makes "spend more, get a better model" a plannable investment rather than a gamble — the single fact underneath the entire build-out.

The loss the whole run minimizes

The cross-entropy of a single prediction — the quantity, summed over trillions of tokens, that *is* pretraining.

loss.py

```
import numpy as np

def softmax(z):
    e = np.exp(np.array(z, float) - max(z)); return e / e.sum()

def cross_entropy(logits, true_token):
    p = softmax(logits)
    return -np.log(p[true_token])    # surprise at the real next token

# model's logits over a tiny vocab; the true next token is index 0
logits = [6.0, 1.5, 1.8, 1.2]
print(round(cross_entropy(logits, 0), 3))    # 0.034 - confident and correct:
print(round(cross_entropy(logits, 2), 3))    # 4.234 - it bet wrong: large loss

# training = average this over D tokens and take a gradient step, ~D/batch time
```

The hundred-million-dollar single event

THE RUN → MONEY

Pretraining is the most concentrated cost in all of AI — a single training run for a frontier model costs **tens to hundreds of millions of dollars** in compute, burned over weeks or months on a cluster of tens of thousands of GPUs that exists largely for this purpose. It is the sharpest spike of the build-out's capital expenditure, and it happens before the model has earned a cent.

The economics only work because the result is an **asset**: one expensive run produces a base model that is then served to hundreds of millions of users, amortizing the cost across billions of cheap-by-comparison [inference](#) calls. The scaling law is what makes the bet rational — spend predictably more on the run, get a predictably better asset.

So pretraining is the training half of the two clocks: the enormous, upfront, capitalized cost that the [Circuit](#) asks whether inference revenue will ever repay. Every new frontier run raises the stakes of that question.

The primary sources

Brown et al. (2020) — GPT-3 · pretraining at scale and few-shot learning.

Hoffmann et al. (2022) — Chinchilla · how to spend a compute budget between parameters and data.

Wei et al. (2022) — Emergent Abilities of LLMs · capabilities that appear with scale.

Gao et al. (2020) — The Pile · what a pretraining corpus actually looks like.

Cite this chapter: Divergent Compute, "How models are pretrained", First Principles, 2026.
divergentcompute.com/first-principles-pretraining · v1.0 · CC-BY.

The math the consensus is skipping — unlocked for life.

Join as a Founding Reader for free lifetime access to the interactive instruments, briefs, and data boards — before the gate closes.

500 LIFETIME FOUNDING SEATS · CLOSING AUG 25, 2026

Institutional email

Secure access

Fine-tuning & RLHF

A pretrained base model knows language but has no manners — it just continues text. Fine-tuning and RLHF are the two phases that turn that raw predictor into a model that **follows instructions, stays helpful, and refuses harm.**

Read at your depth: [01 The answer](#) · [02 Intuition](#) · [03 Mechanics](#) · [04 The math](#) · [05 The code](#)
· [06 The economics](#) · [07 Sources](#)

Teaching a text-predictor to behave like an assistant

Out of pretraining you get a **base model**: brilliant at predicting text, useless as an assistant. Ask it "How do I reset my password?" and it might continue with a list of *more questions* — because that's what such text looks like on the web. It has the knowledge; it lacks the job description.

Two phases fix that. **Supervised fine-tuning (SFT)** shows it thousands of high-quality instruction→response examples, teaching it the *shape* of being helpful. Then **RLHF** — reinforcement learning from human feedback — has people rank competing answers, distills those preferences into a reward signal, and nudges the model toward replies humans actually prefer: clearer, more honest, and safe.

Pick a prompt and walk it through the three stages. The knowledge never changes — only the behaviour:

BASE → SFT → RLHF — SAME PROMPT, EVOLVING BEHAVIOUR

Illustrative responses showing what each alignment phase adds.

How do I reset my password?

Write a haiku about the sea.

How do I pick a lock I don't own?

Prompt: **“How do I reset my password?”**

Base model

+ SFT

+ RLHF

How do I change my email address? How do I update my billing? How do I delete my account? How do I contact support about my login...

02 MECHANICS

Two phases, two kinds of signal

- **SFT (supervised fine-tuning)**. Continue training the base model — same next-token objective — but now on a curated set of `instruction → ideal response` pairs written by humans. The model learns the assistant format: when it sees a question, produce an answer, not more questions. Cheap relative to pretraining, but the data is hand-written and expensive per token.
- **The reward model**. Humans are shown two responses to the same prompt and pick the better one. A separate model is trained on thousands of these comparisons to output a scalar **reward** — a learned proxy for "what humans prefer." Ranking is far easier and more reliable for humans than writing perfect answers.
- **RLHF (policy optimization)**. The model is then optimized — typically with PPO, or more recently the simpler **DPO** — to produce responses the reward model scores highly, while a **KL penalty** keeps it from drifting too far from the sensible SFT model and "gaming" the reward.
- **The result**. Helpfulness, honesty, and harmlessness — the behaviour, not the knowledge. This is "alignment," and it's the difference between a research artifact and a product millions will pay for.

One subtlety worth holding onto: RLHF optimizes for *what raters prefer*, which is a proxy for what's actually good. Push too hard and you get models that are confidently agreeable rather than correct — the central open problem the field calls alignment.

Preferences become a reward

The reward model turns "humans prefer A over B" into numbers via the **Bradley–Terry** model. If the reward model scores responses $r(A)$ and $r(B)$, the modeled probability that A is preferred is:

$$P(A \succ B) = \sigma(r(A) - r(B)) = \frac{1}{1 + e^{-(r(A) - r(B))}}$$

It's trained to maximize that probability on the human-labeled pairs — i.e. minimize $-\log \sigma(r(A) - r(B))$ whenever a human chose A . The policy π is then optimized to earn reward while staying near the reference (SFT) model π_{ref} :

$$\max_{\pi} \mathbb{E}_{y \sim \pi} [r(x, y)] - \beta \text{KL}(\pi \parallel \pi_{\text{ref}})$$

The first term pulls toward what humans like; the βKL term is the leash that stops it from degenerating into reward-hacking gibberish. Tuning β is the whole art — too loose and the model games the reward, too tight and RLHF does nothing.

The preference loss, in nine lines

How a single human comparison becomes a gradient on the reward model.

Runnable.

reward.py

```
import numpy as np

def sigmoid(x): return 1.0 / (1.0 + np.exp(-x))

# reward model's current scores for two candidate responses
r_chosen, r_rejected = 2.3, 0.8      # a human preferred the first one

p_agree = sigmoid(r_chosen - r_rejected)  # model's prob the human's pick is
loss    = -np.log(p_agree)               # Bradley-Terry preference loss

print(round(p_agree, 3))  # 0.818  - model mostly agrees with the human
print(round(loss, 3))    # 0.201  - small loss; a confident wrong pick costs

# training nudges r_chosen up and r_rejected down until the model ranks like p
```

The cheap phase that creates all the value

ALIGNMENT → MONEY

Fine-tuning is a rounding error next to [pretraining](#) in compute — but it is where a model becomes **sellable**. Nobody subscribes to a base model; people pay for the aligned assistant. The hundred-million-dollar pretraining asset only earns revenue after this comparatively cheap polishing step gives it manners.

The cost here is a different shape: not GPUs but **people**. SFT demonstrations and preference labels are written and ranked by human annotators, and high-quality human feedback at scale is a real, recurring line item — increasingly the scarce input as raw compute becomes plentiful. "Data" stops meaning scraped web text and starts meaning curated human judgment.

So alignment is the conversion step in the [Circuit](#): the move that turns a capitalized training asset into a product with paying users. It's also where the quality questions live — whether the helpful, confident answers people are paying for are actually *right* is the thing our research desk keeps testing.

The primary sources

Ouyang et al. (2022) — InstructGPT · the SFT + RLHF recipe behind ChatGPT.

Bai et al. (2022) — Constitutional AI · using AI feedback to reduce the human-labeling burden.

Rafailov et al. (2023) — Direct Preference Optimization (DPO) · skipping the separate reward model.

Christiano et al. (2017) — Deep RL from Human Preferences · the original idea, before LLMs.

Cite this chapter: Divergent Compute, "Fine-tuning & RLHF", First Principles, 2026.
divergentcompute.com/first-principles-finetuning · v1.0 · CC-BY.

The math the consensus is skipping — unlocked for life.

Join as a Founding Reader for free lifetime access to the interactive instruments, briefs, and data boards — before the gate closes.

500 LIFETIME FOUNDING SEATS · CLOSES AUG 25, 2026

Institutional email

Secure access

Differences between LLMs

GPT, Claude, Gemini, Llama, Mistral, DeepSeek — almost all of them are the **same transformer recipe** you just learned. What actually separates them lives on a handful of axes: scale, data, architecture, alignment, openness, and context.

Read at your depth: [01 The answer](#) · [02 Intuition](#) · [03 Mechanics](#) · [04 The math](#) · [05 The code](#) · [06 The economics](#) · [07 Sources](#)

One recipe, a few dials

It's tempting to imagine each lab guards some secret architecture. Mostly they don't. Strip the branding away and you find the same decoder-only transformer, trained with the same next-token objective and aligned the same way. The famous models differ by *degree* on a small set of dials, not by kind.

That's genuinely clarifying: once you understand one model, you understand them all, and the marketing resolves into six measurable choices. Click each axis to see what it controls — and where real, publicly documented models actually sit:

THE SIX AXES OF DIFFERENCE — EXPLORE

Only publicly disclosed facts. Frontier closed-model sizes are undisclosed and shown as such.

- Scale
- Data
- Architecture
- Alignment
- Openness
- Context

Scale — parameter count

How many weights the model has. More capacity, but only useful if matched with enough data. Closed frontier models no longer disclose this.

| | | |
|--------------|---------------------------------------|------------------------|
| small | Mistral 7B | Llama-3 8B |
| mid | Llama-3 70B | Mixtral 8x7B |
| large (open) | Llama-3.1 405B | DeepSeek-V3 671B (MoE) |
| historic | GPT-3 175B (2020) | |
| frontier | GPT-4 / Claude / Gemini — undisclosed | |

Note · Sizes shown are publicly disclosed by the labs. Current frontier closed models do not publish parameter counts.

What each dial actually changes

- **Scale.** More parameters = more capacity, but only if matched with enough data. A bigger model isn't automatically better — an undertrained giant loses to a well-fed smaller one.
- **Data.** The mix, quality, and quantity of training tokens. This is the least visible and arguably most decisive dial — and the hardest to copy, since each lab's data pipeline is proprietary.
- **Architecture tweaks.** The core is shared, but labs vary attention (grouped-query, sliding-window), use **mixture-of-experts** to grow parameters without growing per-token cost, and extend context length. Differences of efficiency, not of kind.
- **Alignment.** The post-training — SFT data, RLHF quality, system behaviour — is where a model's "personality," refusal style, and reliability come from. Two models with identical pretraining can feel completely different after this step.
- **Openness.** Open-weight (Llama, Mistral, Qwen, DeepSeek — you can download and run them) vs closed (GPT, Claude, Gemini — API only). This is a business and safety choice, not a capability one.
- **Context window.** How much the model can attend to at once — from 8K to a million-plus tokens. A real capability difference for long-document work, and a real cost difference.

Benchmarks try to collapse all of this into a single leaderboard number. Treat those with suspicion — they're easily gamed and rarely capture the axis you actually care about for your task.

Why a smaller model can win

Training compute is set by parameters N and tokens D :

$$C \approx 6 N D$$

The **Chinchilla** result says that for a fixed compute budget C , loss is minimized when parameters and data grow *together* — both scale roughly as the square root of compute:

$$N_{\text{opt}} \propto C^{0.5}, \quad D_{\text{opt}} \propto C^{0.5} \quad \Rightarrow \quad D_{\text{opt}} \approx 20 N_{\text{opt}}$$

The rule of thumb that fell out: about **20 tokens per parameter** for compute-optimal training. This is why model size alone tells you little. A 405B model trained on too few tokens is *undertrained*; an 8B model trained on 15 trillion tokens is pushed far past compute-optimal on purpose — because the lab is optimizing not for training cost but for cheap inference later. Different objective, different point on the curve.

The giant a small model out-trained

Compute $6ND$ for four models with publicly disclosed sizes and token counts. Note the surprise in the output.

compute.py

```
# training compute C = 6 * N(params) * D(tokens) – all figures publicly disclosed
models = [
    ("GPT-3 (2020)", 175e9, 300e9), # 175B params, 300B tokens
    ("Llama-3 8B", 8e9, 15e12), # 8B params, 15T tokens
    ("Llama-3 70B", 70e9, 15e12),
    ("Llama-3.1 405B", 405e9, 15e12),
]
for name, N, D in models:
    C = 6 * N * D
    print(f"{name:16s} N={N:.0e} D={D:.0e} C={C:.2e} FLOPs")
# GPT-3 (2020)      N=2e+11 D=3e+11 C=3.15e+23 FLOPs
# Llama-3 8B       N=8e+09 D=2e+13 C=7.20e+23 FLOPs  <- MORE compute than GPT-3
# Llama-3 70B      N=7e+10 D=2e+13 C=6.30e+24 FLOPs    at 1/22 the size
# Llama-3.1 405B  N=4e+11 D=2e+13 C=3.65e+25 FLOPs
```

Llama-3 8B is **22x smaller** than GPT-3 yet consumed **~2.3x more training compute**, because it was fed 50x the data. Parameter count is a label on the box, not a measure of what went into it.

The market the dials create

DIFFERENCES → MONEY

These axes aren't just technical — they carve the market. **Open-weight** models commoditize the bottom: anyone can run Llama or Mistral on their own hardware, so the price of "good enough" intelligence collapses toward the cost of electricity. **Closed frontier** models defend a premium at the top, charging per token for the best available quality. The whole industry is a tug-of-war between those two forces.

The compute-optimal math is why: labs now deliberately *overtrain* small models, spending more upfront so that [inference](#) — the part you pay for forever — is cheap. A small, heavily-trained model that runs on one GPU can undercut a frontier API on price for most everyday tasks. That pressure on margins, even as capability rises, is exactly the tension the [Circuit](#) tracks.

So "which model is best" is the wrong question economically. The right one is which point on these six dials fits your task and your budget — and whether the frontier premium survives as the open tier keeps catching up.

The primary sources

Llama 3 Herd of Models (2024) · disclosed sizes, token counts, and design choices.

Hoffmann et al. (2022) — Chinchilla · the compute-optimal ~20-tokens-per-parameter rule.

Mixtral of Experts (2024) · mixture-of-experts: more parameters, same per-token cost.

DeepSeek-V3 (2024) · an open MoE frontier model, with disclosed training cost.

Cite this chapter: Divergent Compute, "Differences between LLMs", First Principles, 2026.
divergentcompute.com/first-principles-model-differences · v1.0 · CC-BY.

The math the consensus is skipping — unlocked for life.

Join as a Founding Reader for free lifetime access to the interactive instruments, briefs, and data boards — before the gate closes.

500 LIFETIME FOUNDING SEATS · CLOSING AUG 25, 2026

Secure access

Multimodal models

A model that "sees" doesn't need a new brain. An image is just chopped into patches, each patch turned into a token in the same vector space as words — and fed into the **same transformer**. Everything becomes tokens.

Read at your depth: [01 The answer](#) · [02 Intuition](#) · [03 Mechanics](#) · [04 The math](#) · [05 The code](#)
· [06 The economics](#) · [07 Sources](#)

A picture is just more tokens

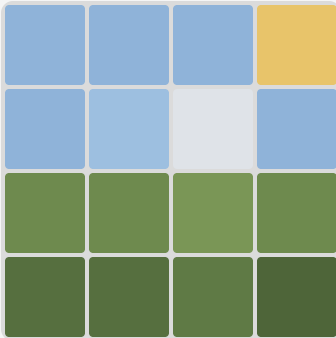
The deep trick of multimodal models is that there's no trick. The transformer doesn't care whether a token came from a word or a corner of a photo — it just attends over a sequence of vectors. So to make a model "see," you convert an image into a sequence of vectors and splice them in next to the text tokens.

The conversion is mechanical: cut the image into a grid of small **patches**, flatten each patch, and run it through a linear layer that projects it into the model's embedding dimension — exactly where word tokens live. Now text and image share one sequence, and attention can relate "the cat" in the prompt to the patch of the photo where the cat is.

Click any patch of this little image and watch it light up as a token sitting right alongside the words:

ONE IMAGE → PATCHES → TOKENS IN THE SAME SEQUENCE

A stylized 4x4 image. Click a patch (or a token) — each patch is exactly one token.



INPUT SEQUENCE TO THE TRANSFORMER

Describe this image : + [blue] [blue] [blue] [yellow] [blue]

[blue] [light gray] [blue] [green] [green] [green] [green] [green] [green] [green]

Click a patch to see it become a token.

↓ patchify & project

From pixels to a shared sequence

- **Patchify.** Split the image (say 224×224) into fixed patches (say 16×16). That yields a grid of patches — for those numbers, $14 \times 14 = 196$ of them.
- **Embed.** Flatten each patch's pixels and pass them through one linear layer that maps them to the model's embedding dimension d . Add a position embedding so the model knows where each patch sat. This little vision front-end is a **Vision Transformer (ViT)**.
- **Concatenate.** Place the image tokens into the sequence with the text tokens: `[text... , img_1, img_2, ... , img_196]`. From here the main transformer is unchanged — it just attends over a longer, mixed sequence.
- **Train the bridge.** The projection (and often the whole stack) is trained on image–text pairs so the visual tokens land in a space the language model already understands — the lineage of **CLIP** and models like LLaVA.
- **Other modalities, same idea.** Audio becomes spectrogram patches; video becomes frames-worth of patches across time. Generation runs the other way — producing image or audio tokens, often via a diffusion or codec decoder.

So "multimodal" is less a new architecture than a set of **adapters** that turn every kind of input into tokens the one transformer can read.

Patches, counted and projected

An image of height H and width W , cut into square patches of side P , yields a number of patch tokens:

$$N_{\text{patches}} = \frac{H}{P} \times \frac{W}{P} = \frac{H W}{P^2}$$

Each patch is a block of $P \times P \times C$ pixel values (C = colour channels). Flattened to a vector $x_p \in \mathbb{R}^{P^2 C}$, it's projected to the embedding dimension d by a learned matrix E and given a position embedding:

$$z_p = E x_p + \text{pos}_p, \quad E \in \mathbb{R}^{d \times P^2 C}$$

The full input is then the concatenation $[z_1^{\text{text}}, \dots, z_m^{\text{text}}, z_1^{\text{img}}, \dots, z_N^{\text{img}}]$ — one sequence of d -dimensional vectors. Because attention cost grows with the *square* of sequence length, those N image tokens are not free.

How many tokens is a picture?

Patch counts for real input sizes — the number that lands in your context window per image.

patches.py

```
def patch_tokens(H, W, P):
    return (H // P) * (W // P)          # N = (H*W) / P^2

for H, W, P in [(224, 224, 16), (336, 336, 14), (512, 512, 16)]:
    print(f"{H}x{W}, patch {P}: {patch_tokens(H, W, P)} tokens")
# 224x224, patch 16: 196 tokens
# 336x336, patch 14: 576 tokens
# 512x512, patch 16: 1024 tokens    <- one image ≈ a page of text, in tokens

# a 10-word text prompt is ~13 tokens; a single image can be 50x larger
```

Why a photo costs like a page

PIXELS → MONEY

Because images are billed as tokens, vision is expensive. A single picture can be **hundreds to over a thousand tokens** — a short text prompt is a dozen. So sending one photo can cost more than a long paragraph, and the [quadratic attention bill](#) means a few high-resolution images can dominate a request's entire compute.

Video is the extreme case: it's images stacked through time, so a few seconds of footage can be tens of thousands of tokens. This is why "just feed it the whole video" is an economic, not just technical, problem — and why providers downsample frames aggressively. Multimodality is one of the largest forces inflating the [inference](#) demand the build-out is racing to supply.

The strategic read for the [Circuit](#): every new modality multiplies the tokens per interaction, which multiplies the compute per user, which raises the revenue each interaction must eventually justify. Seeing and hearing are powerful — and they move the break-even line further out.

The primary sources

Dosovitskiy et al. (2020) — An Image is Worth 16×16 Words (ViT) · images as patch tokens.

Radford et al. (2021) — CLIP · aligning image and text in one embedding space.

Alayrac et al. (2022) — Flamingo · bridging a vision encoder into a language model.

Liu et al. (2023) — LLaVA · a simple, open visual-instruction recipe.

Cite this chapter: Divergent Compute, "Multimodal models", First Principles, 2026.
divergentcompute.com/first-principles-multimodal · v1.0 · CC-BY.

The math the consensus is skipping — unlocked for life.

Join as a Founding Reader for free lifetime access to the interactive instruments, briefs, and data boards — before the gate closes.

500 LIFETIME FOUNDING SEATS · CLOSING AUG 25, 2026

Secure access

Quantization & distillation

A frontier model is enormous. Two techniques shrink it to run cheaply: **quantization** stores each weight in fewer bits, and **distillation** trains a small "student" to copy a big "teacher" — keeping most of the quality at a fraction of the size.

Read at your depth: [01 The answer](#) · [02 Intuition](#) · [03 Mechanics](#) · [04 The math](#) · [05 The code](#)
· [06 The economics](#) · [07 Sources](#)

Same model, fewer bits — and a smaller copycat

From the [parameters](#) chapter you know the rule: memory = parameter count × bytes per parameter. The second factor is the lever. By default each weight is a 16-bit number, but most of that precision is wasted — the model works almost as well if each weight is squeezed into 8 bits, or even 4. That's **quantization**, and it cuts memory and cost in direct proportion.

Distillation attacks the first factor. Instead of compressing one model, you train a much smaller one to imitate the big model's outputs — learning from the teacher's full probability distribution, not just the right answer. The student ends up far smaller while inheriting much of the teacher's skill.

Quantization is the bigger everyday lever. Pick a model and drag the precision down — watch the hardware bill collapse:

QUANTIZATION CALCULATOR — DRAG THE PRECISION

Memory = parameters × (bits / 8). GPUs assume 80 GB each with ~20% runtime overhead.

7B model 70B model 175B model

140GB
WEIGHTS IN MEMORY

3
80GB GPUS NEEDED

baseline
VS 16-BIT

precision 16-bit

4-bit 8-bit 16-bit 32-bit

16-bit · the default. What "weights × 2 bytes" assumes.

Compressing the weights, and the knowledge

- **Quantization.** Map each high-precision weight to a low-bit integer by storing a shared scale per group of weights: $w \approx \text{scale} \times \text{round}(w / \text{scale})$. 16→8 bits roughly halves memory for almost no quality loss; 4-bit (with careful methods like GPTQ or AWQ) quarters it with a small, often acceptable hit. The weights are the same model — just stored coarsely.
- **Why it works.** Neural networks are remarkably robust to noise in their weights. The signal lives in the overall pattern, not the 12th decimal place, so throwing away low-order bits costs surprisingly little.
- **Distillation.** Run the big **teacher** on lots of inputs and record its full output distribution (the "soft labels"). Train a small **student** to match those distributions. The soft targets carry far more information than a single correct token — the teacher reveals *how* confident it is across all options — so the student learns faster and smaller.
- **Used together.** Modern small models are often distilled from a big teacher and *then* quantized — the two compressions compound, which is how genuinely capable models end up running on a laptop or phone.

Both are lossy. The craft is spending the loss where it doesn't matter — and measuring honestly whether the smaller model still does your task, rather than trusting a benchmark.

Bits, scales, and soft targets

Memory for the weights is linear in the bit-width b :

$$\text{memory} = N \times \frac{b}{8} \text{ bytes}$$

Quantization to b bits picks a scale s for a group of weights and rounds each to the nearest representable level, then dequantizes on use:

$$w_q = \text{round}\left(\frac{w}{s}\right), \quad \hat{w} = s w_q \approx w$$

Distillation trains the student distribution p_S to match the teacher's softened distribution p_T (temperature T flattens both), minimizing the KL divergence:

$$\mathcal{L}_{\text{distill}} = \text{KL}(p_T \parallel p_S) = \sum_i p_T(i) \log \frac{p_T(i)}{p_S(i)}$$

The teacher's *soft* probabilities — "70% cat, 25% dog, 5% fox" — teach the student the structure a hard label ("cat") never could. That extra signal is why a small student can punch so far above its size.

The bill, bit by bit

The same 70B model at four precisions — memory and GPUs. This is the calculator above, in seven lines.

quantize.py

```
import math

def mem_gb(N, bits): return N * bits / 8 / 1e9
def gpus(m, cap=80, overhead=1.2): return math.ceil(m * overhead / cap)

N = 70e9 # a 70-billion-parameter model
for bits in [32, 16, 8, 4]:
    m = mem_gb(N, bits)
    print(f"{bits}>2}-bit: {m:6.1f} GB -> {gpus(m)} GPU(s)")
# 32-bit: 280.0 GB -> 5 GPU(s)
# 16-bit: 140.0 GB -> 3 GPU(s)
# 8-bit: 70.0 GB -> 2 GPU(s)
# 4-bit: 35.0 GB -> 1 GPU(s) <- a 70B model on a single GPU
```

The deflation underneath everything

COMPRESSION → MONEY

Quantization and distillation are the **deflationary force** in AI. Every bit you drop is a proportional cut in memory, hardware, and energy per token served — and distillation moves a capability from an expensive frontier model into one a fraction of the size. Together they are why the price of a given level of intelligence keeps falling fast even as the frontier rises.

This is the optimistic half of the Circuit's ledger. The build-out spends ever more on the frontier, but compression relentlessly drives down the cost of serving everything below it — pushing models onto single GPUs, laptops, and phones, where the marginal cost approaches electricity. The open tier rides this hardest.

It also sharpens the central tension. If a distilled, quantized open model captures most of the value at a tenth of the cost, what exactly is the frontier premium being paid for — and for how long? That question, asked with real numbers, is the whole point of this think tank.

You now know what a model *is* — and what it costs

From a base model learning language by guessing the next word, through alignment, the real differences between models, vision, and the compression that makes any of it affordable — you've followed the model from raw weights to a product with a price tag. **Part I** built the machine; **Part II** made it a model.

Part III opens the machinery that runs it — inference, GPUs, the memory wall, the KV cache, and the data-center cluster: what it actually takes to serve a model to millions at once, and where the real costs live. [See the full curriculum →](#)

07 GOING DEEPER

The primary sources

Hinton et al. (2015) — Distilling the Knowledge in a Neural Network · the soft-target idea.

Dettmers et al. (2022) — LLM.int8() · 8-bit inference for large models with no quality loss.

Frantar et al. (2022) — GPTQ · accurate 4-bit post-training quantization.

Dettmers et al. (2023) — QLoRA · fine-tuning a 4-bit model on a single GPU.

Cite this chapter: Divergent Compute, "Quantization & distillation", First Principles, 2026.
divergentcompute.com/first-principles-quantization · v1.0 · CC-BY.

The math the consensus is skipping — unlocked for life.

Join as a Founding Reader for free lifetime access to the interactive instruments, briefs, and data boards — before the gate closes.

500 LIFETIME FOUNDING SEATS · CLOSES AUG 25, 2026

Institutional email

Secure access

What is inference?

Training builds the model once. **Inference** is running it — the forward passes that turn your prompt into an answer, every single time you hit send. It happens in two very different phases, and it is the cost you pay forever.

Read at your depth: [01 The answer](#) · [02 Intuition](#) · [03 Mechanics](#) · [04 The math](#) · [05 The code](#)
· [06 The economics](#) · [07 Sources](#)

Read the prompt fast; write the answer slow

Inference splits into two phases that feel nothing alike. **Prefill** reads your entire prompt in a single parallel pass — every token processed at once, the GPU running flat out. Then **decode** begins, and the model writes its answer one token at a time: each new token needs its own forward pass, and each pass depends on the one before, so they can't be parallelized.

That asymmetry is the whole story of inference cost. Reading is cheap and parallel; writing is sequential and slow. Click **Run inference** and watch it happen — the prompt lights up all at once, then the answer crawls out token by token:

RUN INFERENCE — PREFILL, THEN DECODE

One parallel prefill pass over the prompt, then one forward pass per generated token.

PROMPT & GENERATED ANSWER

What is inference ?

| | | |
|----------------|---------------|----------------------|
| 0 | 0 | 0 |
| PREFILL PASSES | DECODE PASSES | TOTAL FORWARD PASSES |

Run inference Reset

Two phases, two bottlenecks

- **Prefill (compute-bound).** The whole prompt is pushed through the model in one pass. Because all prompt tokens are available at once, the math is one big parallel matrix multiply — the GPU's compute is the limit, and it's used efficiently. This phase sets the "time to first token."
- **Decode (memory-bound).** Now the model generates. Each token requires loading the model's weights from memory to do one forward pass, producing exactly one token, which is appended and fed back in. The arithmetic per step is tiny relative to the weights moved, so the bottleneck flips from compute to **memory bandwidth** — the GPU mostly waits on memory. This is why generation streams out at a steady, limited pace.
- **The KV cache.** To avoid recomputing the whole prompt at every decode step, the model caches the per-token intermediate values (keys and values) — the KV cache. It's what makes decode merely slow instead of catastrophically slow, and it's a major consumer of GPU memory (Chapter 18).
- **No learning happens.** Inference only runs the forward pass — no gradients, no weight updates. The model is frozen; it's a pure function from tokens to tokens.

So the model you spent \$100M training does the same forward-pass arithmetic on every request for the rest of its life. Making that arithmetic fast and cheap is what all of Part III is about.

Two FLOPs per parameter, per token

One forward pass through a model with N parameters costs about $2N$ floating-point operations per token (a multiply and an add per weight). For a request with a prompt of P tokens and a generated answer of G tokens:

$$\text{prefill} \approx 2N P \quad (1 \text{ parallel pass}), \quad \text{decode} \approx 2N G \quad (G \text{ sequential passes})$$

$$\text{total} \approx 2N (P + G)$$

The FLOP counts can be equal, but the *wall-clock* isn't: prefill's P tokens run together, while decode's G tokens run strictly one after another. The decode phase is also **memory-bound** — its limiter isn't FLOPs but the bytes of weights streamed per token, roughly $2N \cdot$ (bytes per parameter) moved from memory each step. That ratio, arithmetic intensity, is why decode leaves a GPU's compute mostly idle — and why fewer bits and batching matter so much.

The cost of one request

Prefill vs decode FLOPs for a 70B model answering with a 500-token prompt and 500-token reply.

inference.py

```
def inference_flops(N, prompt_len, gen_len):
    per_token = 2 * N          # ~2 FLOPs per parameter, per token
    prefill = per_token * prompt_len  # all prompt tokens, ONE parallel pass
    decode = per_token * gen_len      # gen tokens, that many SEQUENTIAL passes
    return prefill, decode

N = 70e9
pf, dc = inference_flops(N, 500, 500)
print(f"prefill: {pf:.2e} FLOPs (1 parallel pass, 500 prompt tokens)")
print(f"decode: {dc:.2e} FLOPs (500 sequential passes)")
print(f"total: {pf+dc:.2e} FLOPs per request")
# prefill: 7.00e+13 FLOPs (1 parallel pass, 500 prompt tokens)
# decode: 7.00e+13 FLOPs (500 sequential passes)
# total: 1.40e+14 FLOPs per request <- and decode's are one-at-a-time
```

The cost that never stops

INFERENCE → MONEY

Training is a one-time capital event. **Inference is the bill that arrives forever** — every request, from every user, runs the full forward pass again. Across a model's life, serving it almost always costs far more in total than training it did, which is why the data centers being built are sized for *inference* demand, not just training runs.

And the expensive half is decode. Because it's sequential and memory-bound, a single user's generation barely uses a GPU's compute — so providers pack many users' requests together (batching, Chapter 18) to fill the silicon. The entire discipline of inference engineering exists to claw back the efficiency that the one-token-at-a-time nature of decode throws away.

This is the meter the Circuit watches most closely: every token decoded is a real, recurring cost, and the question is whether the revenue per token clears it. Training built the asset; inference is where the money is actually spent — and, one hopes, made.

The primary sources

Pope et al. (2022) — Efficiently Scaling Transformer Inference · prefill vs decode, the cost model.

Kwon et al. (2023) — PagedAttention / vLLM · the modern serving engine.

Kaplan et al. (2020) — Scaling Laws · the $2N$ FLOPs-per-token accounting.

How to Scale Your Model (Google DeepMind) · a deep, free reference on inference arithmetic.

Cite this chapter: Divergent Compute, "What is inference?", First Principles, 2026.
divergentcompute.com/first-principles-inference · v1.0 · CC-BY.

The math the consensus is skipping — unlocked for life.

Join as a Founding Reader for free lifetime access to the interactive instruments, briefs, and data boards — before the gate closes.

500 LIFETIME FOUNDING SEATS · CLOSES AUG 25, 2026

Secure access

Latency, throughput, tokens/sec

How fast is an AI model? The honest answer is a formula. Because decode is memory-bound, a model's speed is set by how fast it can stream its weights out of memory — **tokens per second \approx bandwidth \div model bytes**.

Read at your depth: [01 The answer](#) · [02 Intuition](#) · [03 Mechanics](#) · [04 The math](#) · [05 The code](#)
· [06 The economics](#) · [07 Sources](#)

Speed is bandwidth divided by size

You felt this in the last chapter: decode writes one token per forward pass, and each pass has to read the *entire* model out of memory. So the speed limit isn't how fast the GPU can compute — it's how fast it can **move the weights**. A bigger model means more bytes to stream per token, which means fewer tokens per second. It's almost that simple.

Two numbers matter to a user: **time to first token** (how long until the answer starts — set by prefill) and **tokens per second** (how fast it then types — set by decode). The calculator below computes the second from first principles. Pick a model size and drag the precision; watch a 70B model crawl, and watch [quantization](#) nearly quadruple it:

TOKENS/SEC CALCULATOR — THE MEMORY-BOUND SPEED LIMIT

Single stream on one GPU with ~3.35 TB/s memory bandwidth (H100-class). tok/s \approx bandwidth \div (params \times bytes).

7B model
70B model
175B model

23.9 tok/s

DECODE TOKENS / SEC

41.8 s

TIME FOR A 1,000-TOKEN REPLY

140 GB

BYTES STREAMED / TOKEN

precision 16-bit

16-bit 8-bit 4-bit

24 tok/s · Readable pace — roughly human reading speed. Quantizing or batching makes the economics work.

Latency, throughput, and the tension between them

- **Time to first token (TTFT).** Set by prefill — how long to process your whole prompt before the first word appears. Grows with prompt length; this is the "thinking..." pause.
- **Inter-token latency / tokens-per-second.** Set by decode. Each token streams the model's weights from memory, so per-user speed \approx memory bandwidth \div model bytes. This is the steady typing rate you watch.
- **End-to-end latency.** Simply $\text{TTFT} + (\text{tokens} \times \text{per-token time})$. A long answer from a big model is slow on both counts.
- **Throughput.** The *system's* total tokens/sec across all users at once. Here's the twist: because one user's decode barely uses the GPU's compute, you can run many users' requests in the same forward pass — **batching** — and total throughput climbs almost for free, even though each individual user's speed is unchanged or slightly slower.
- **The tension.** Latency is one user's experience; throughput is the system's efficiency. Bigger batches raise throughput (and lower cost per token) but can raise latency. Tuning that trade-off is the central job of an inference team.

So "tokens per second" is two metrics wearing one name: a *per-user* speed bounded by bandwidth, and an *aggregate* throughput bounded by compute once the batch is full. Quantization helps the first; batching helps the second.

The roofline of decode

Each decode step streams roughly the whole model — N parameters at $b/8$ bytes each — from memory. With memory bandwidth BW (bytes/sec), the time per token and the rate are:

$$t_{\text{token}} \approx \frac{N \cdot (b/8)}{\text{BW}}, \quad \text{tokens/sec} \approx \frac{\text{BW}}{N \cdot (b/8)}$$

End-to-end latency for a prompt of P and an answer of G tokens:

$$\text{latency} \approx \underbrace{t_{\text{prefill}}(P)}_{\text{TTFT}} + G \cdot t_{\text{token}}$$

The roofline idea: decode lives on the memory-bandwidth slope, so halving the bytes (quantize) roughly **doubles** the rate. Batching B requests reuses the same weight-load across all of them, so aggregate throughput rises toward $B\times$ — until the work becomes compute-bound and hits the GPU's FLOPs ceiling instead. Inference engineering is the art of climbing that roofline.

Speed, from the bandwidth up

Single-stream decode rate for three model sizes at three precisions, on an H100-class GPU.

tokens_per_sec.py

```
BW = 3.35e12    # H100 HBM3 memory bandwidth, ~3.35 TB/s

def toks_per_sec(N, bits):
    bytes_streamed = N * (bits / 8)    # ~whole model moved per token
    return BW / bytes_streamed

for N, nm in [(7e9,"7B"), (70e9,"70B"), (175e9,"175B")]:
    rates = " ".join(f"{b}-bit {toks_per_sec(N,b):5.1f}" for b in [16,8,4])
    print(f"{nm:5s}: {rates} tok/s")
# 7B   : 16-bit 239.3   8-bit 478.6   4-bit 957.1   tok/s
# 70B  : 16-bit  23.9   8-bit  47.9   4-bit  95.7   tok/s
# 175B : 16-bit   9.6   8-bit  19.1   4-bit  38.3   tok/s

print(f"70B 16-bit, 1000-token reply: {1000/toks_per_sec(70e9,16):.1f} s") #
```

Throughput is the unit of margin

TOKENS/SEC → MONEY

A GPU costs the same per hour whether it serves one user or a hundred. So the number that decides whether inference makes money is **tokens per second per GPU** — and therefore tokens per dollar. Everything in this chapter is a lever on that number: quantization raises per-user speed, batching raises aggregate throughput, and both cut the cost of every token served.

This is why providers obsess over batching and why they price output tokens the way they do: a memory-bound GPU running a single user is mostly idle silicon being paid for. Filling it is the difference between a gross margin and a loss. The same model can be a profitable product or a money pit depending entirely on how well its operator climbs this roofline.

For the [Circuit](#), tokens/sec/dollar is the efficiency term in the whole equation: as it improves — through better hardware, lower precision, and smarter serving — the cost side of the ledger falls, buying the build-out more time for revenue to catch up. It is the most quietly important number in AI economics.

The primary sources

Pope et al. (2022) — Efficiently Scaling Transformer Inference · latency vs throughput trade-offs.

Databricks — LLM Inference Performance Engineering · TTFT, tokens/sec, batching in practice.

NVIDIA — GPU Performance Background (roofline) · memory-bound vs compute-bound.

MLPerf Inference · the standard latency/throughput benchmark.

Cite this chapter: Divergent Compute, "Latency, throughput, tokens/sec", First Principles, 2026. divergentcompute.com/first-principles-latency · v1.0 · CC-BY.

The math the consensus is skipping — unlocked for life.

Join as a Founding Reader for free lifetime access to the interactive instruments, briefs, and data boards — before the gate closes.

500 LIFETIME FOUNDING SEATS · CLOSES AUG 25, 2026

Secure access

Why AI needs GPUs

A neural network is, underneath, one operation repeated trillions of times: **matrix multiplication**. GPUs win because they do that one thing with thousands of cores at once — while a CPU, built for a few fast sequential tasks, does it a handful at a time.

Read at your depth: [01 The answer](#) · [02 Intuition](#) · [03 Mechanics](#) · [04 The math](#) · [05 The code](#)
· [06 The economics](#) · [07 Sources](#)

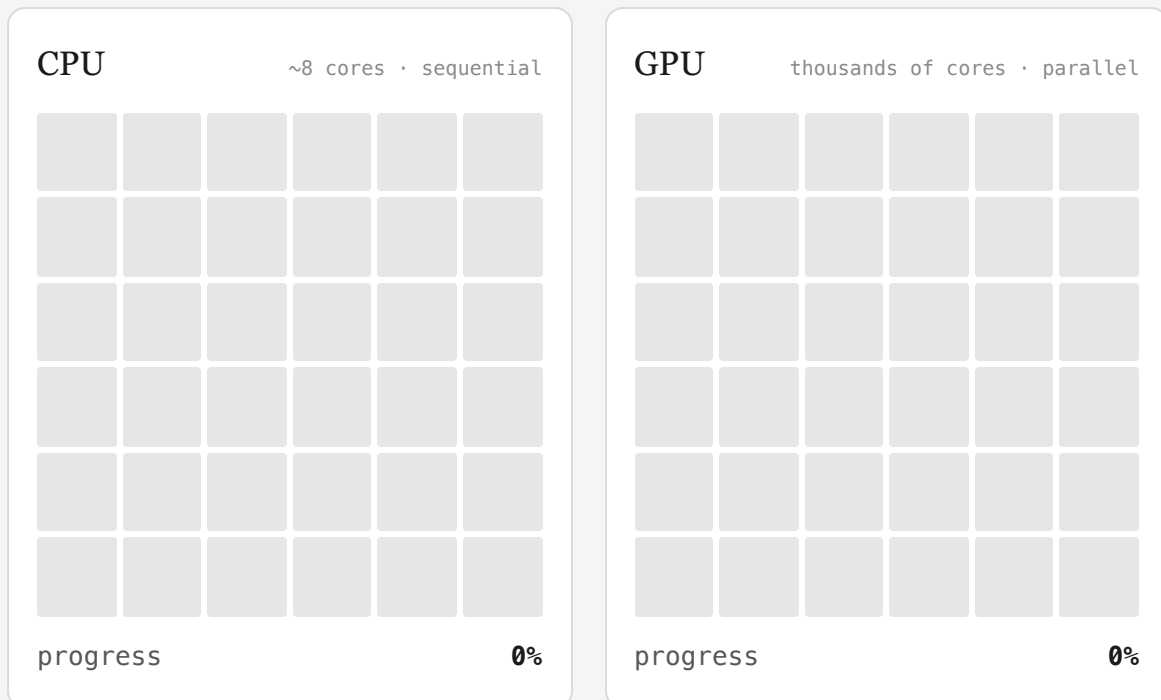
The same sum, a million times over

Strip a transformer down and almost all of its work is multiplying big grids of numbers — every attention step and every weight layer is a matrix multiply. And a matrix multiply is **embarrassingly parallel**: each output number is an independent dot product that doesn't depend on the others, so in principle you could compute all of them at the same time.

A CPU has a few very powerful cores tuned for fast, branchy, sequential work. A GPU has *thousands* of simpler cores built to run the exact same arithmetic on many numbers at once — plus huge memory bandwidth to feed them. For AI's one repeated operation, that's a perfect match. Race them on the same matrix multiply:

CPU VS GPU — THE SAME MATRIX MULTIPLY

Illustrative. Each tile is a chunk of output; cores fill tiles in parallel. The animation shows the *shape* of the gap — the real number is in the readout.



Built for throughput, not latency

- **Thousands of cores.** A modern GPU has tens of thousands of arithmetic units running the same instruction across different data at once (SIMT). For one operation repeated over millions of numbers, that's ideal; for tangled, branchy logic, it's wasted.
- **Tensor cores.** GPUs now include dedicated units that do a small *matrix multiply-accumulate* in a single step — hardware built specifically for the one operation neural networks need. This is most of where the headline FLOP/s comes from.
- **Memory bandwidth.** All those cores are useless if you can't feed them. GPUs pair the compute with very fast high-bandwidth memory (HBM) — the subject of the next chapter — so weights and activations stream in fast enough to keep the cores busy.
- **The trade-off.** A CPU minimizes the time for *one* task (latency); a GPU maximizes the total work done across *many* (throughput). AI is almost entirely the second kind of problem, which is why the GPU — originally built to shade millions of pixels in parallel — turned out to be the perfect AI chip by accident.

This is also why the field is a **hardware** story as much as a software one: progress is gated by how many parallel multiply-adds per second per dollar the silicon can deliver.

Why the gap is $\sim 2,000\times$

Multiplying two $d \times d$ matrices produces d^2 outputs, each a dot product of length d (a multiply and an add per term), so the cost is:

$$\text{FLOPs} = 2 d^3$$

Wall-clock time is just that divided by the hardware's throughput. For $d = 8192$ that's about 1.1×10^{12} FLOPs. A CPU delivering ~ 0.5 TFLOP/s of usable fp32 takes seconds; a GPU's tensor cores deliver on the order of $\sim 1,000$ TFLOP/s:

$$t = \frac{2 d^3}{\text{throughput}} \Rightarrow t_{\text{CPU}} \approx 2.2 \text{ s}, \quad t_{\text{GPU}} \approx 1.1 \text{ ms}$$

A roughly $2,000\times$ difference on a single operation — and a model is billions of these. The gap isn't that the GPU's individual cores are faster (they're slower); it's that there are thousands of them doing independent work at once. Parallelism, not clock speed, is the whole game.

The race, in numbers

One large matrix multiply, timed against a CPU's and a GPU's throughput.

cpu_vs_gpu.py

```
def matmul_flops(d):
    return 2 * d**3          # d^2 outputs, each a length-d dot product

d = 8192
flops = matmul_flops(d)

cpu = 0.5e12      # CPU: ~0.5 TFLOP/s usable fp32
gpu = 990e12     # GPU tensor cores: ~990 TFLOP/s bf16 (H100-class)

print(f"{d}x{d} matmul: {flops:.2e} FLOPs")
print(f"CPU: {flops/cpu:.3f} s")
print(f"GPU: {flops/gpu*1000:.3f} ms")
print(f"speedup: {(flops/cpu)/(flops/gpu):.0f}x")
# 8192x8192 matmul: 1.10e+12 FLOPs
# CPU: 2.199 s
# GPU: 1.111 ms
# speedup: 1980x
```

The build-out is a GPU build-out

PARALLELISM → MONEY

Because AI's appetite is for parallel matrix-multiply throughput, the entire build-out resolves to one scarce thing: **GPUs**. The hundreds of billions in [capital expenditure](#) are, concretely, orders for accelerators and the power and buildings to run them. When people say a lab is "compute-constrained," they mean it cannot get enough of these chips.

That scarcity is why a single company, Nvidia, captures so much of the value in AI — it sells the one input everyone needs — and why its data-center revenue became a real-time gauge of the build-out itself. The chip is the bottleneck, the capital line, and the moat all at once.

For the [Circuit](#), this is the supply side of the equation: the cost of intelligence is set by GPU throughput per dollar, and it falls only as fast as the silicon improves. Every [efficiency trick](#) in this section is ultimately about extracting more useful tokens from each very expensive chip.

The primary sources

NVIDIA — GPU Performance Background · cores, throughput, and the roofline.

NVIDIA — Tensor Core architecture · the matrix-multiply-accumulate unit.

Dao et al. (2022) — FlashAttention · why memory movement, not FLOPs, often dominates.

SemiAnalysis · ongoing economic analysis of the AI hardware supply chain.

Cite this chapter: Divergent Compute, "Why AI needs GPUs", First Principles, 2026.
divergentcompute.com/first-principles-gpus · v1.0 · CC-BY.

The math the consensus is skipping — unlocked for life.

Join as a Founding Reader for free lifetime access to the interactive instruments, briefs, and data boards — before the gate closes.

500 LIFETIME FOUNDING SEATS · CLOSING AUG 25, 2026

Institutional email

Secure access

Memory, HBM & the memory wall

Chips can now do arithmetic far faster than they can fetch the numbers to work on. That widening gap is the **memory wall** — and it's why AI accelerators are defined by their **high-bandwidth memory (HBM)** as much as their compute.

Read at your depth: [01 The answer](#) · [02 Intuition](#) · [03 Mechanics](#) · [04 The math](#) · [05 The code](#) · [06 The economics](#) · [07 Sources](#)

Compute outran memory — and never looked back

For decades, the arithmetic units on chips got faster far quicker than the memory feeding them. Compute has roughly tripled each hardware generation while memory bandwidth grew only about 1.6× — so with every generation, the processor spends *more* of its time waiting on data and less doing math. That is the memory wall, and modern AI slammed straight into it: decode is memory-bound precisely because it moves the whole model per token but does little arithmetic with it.

Drag through the generations and watch the two bars diverge — compute racing ahead, bandwidth trailing, and the gap you must somehow bridge growing wider each step:

THE MEMORY WALL — COMPUTE VS BANDWIDTH, PER GENERATION

Illustrative model of the documented trend (Gholami et al. 2024): compute $\times 3.0$ / generation, bandwidth $\times 1.6$ / generation. Bars are relative to generation 0.



The gap: **$\times 1.00$**

02 MECHANICS

What HBM is, and why it's never enough

- **The problem.** A GPU's thousands of cores can only work as fast as memory delivers operands. When the arithmetic is cheap but the data movement is expensive, the cores sit idle — "memory-bound." Most LLM inference lives here.
- **HBM (high-bandwidth memory).** The fix is to stack DRAM dies vertically, right next to the compute, connected by an extremely wide interface. That buys enormous bandwidth — an H100's HBM moves ~3.35 TB/s, versus tens of GB/s for a normal CPU's DRAM. It's the single feature that makes a GPU a viable AI chip.
- **The catches.** HBM is expensive to manufacture, limited in capacity (tens of GB per chip), power-hungry, and made by only a handful of suppliers. So you're always short of it — short of *capacity* (the weights plus the KV cache must fit) and short of *bandwidth* (they must stream fast enough).
- **Why it keeps biting.** Because compute keeps outpacing bandwidth, each new generation makes the imbalance worse, not better. Chip designers respond with bigger caches, lower precision, and clever data movement — but the wall recedes slowly at best.

So the mental model flips: for AI, a chip's memory system is often the real product, and the compute is the part that's easy to oversupply.

The roofline and the ridge point

The **roofline model** says achievable performance is capped by whichever runs out first — compute or bandwidth — given a workload's **arithmetic intensity** I (FLOPs done per byte moved):

$$\text{performance} = \min(\text{peak compute}, I \times \text{bandwidth})$$

The crossover — the **ridge point** — sits at:

$$I^* = \frac{\text{peak compute}}{\text{bandwidth}}$$

A workload is memory-bound when $I < I^*$. As compute grows $\times 3$ per generation and bandwidth only $\times 1.6$, I^* climbs by $\times(3/1.6) \approx 1.9$ each generation — the bar for staying compute-bound keeps rising. LLM decode at batch 1 has an intensity of roughly 1–2 FLOPs/byte, far below a modern GPU's ridge point of several hundred — which is exactly why it wastes most of the chip's compute, and why batching (raising I) is the escape.

The wall, generation by generation

Relative compute, bandwidth, and the widening gap, using the documented per-generation growth rates.

memory_wall.py

```
# Gholami et al. "AI and Memory Wall" (2024): per ~2-year generation,
# peak compute grows ~3x, memory bandwidth ~1.6x.
C, B = 3.0, 1.6

for k in range(6):
    comp, bw = C**k, B**k
    gap = comp / bw          # extra arithmetic-per-byte you must find to stay
    print(f"gen {k} (+{2*k}yr): compute x{comp:6.1f} bandwidth x{bw:5.1f} gap x{gap:5.1f}")
# gen 0 (+0yr): compute x  1.0 bandwidth x  1.0 gap x 1.00
# gen 3 (+6yr): compute x 27.0 bandwidth x  4.1 gap x 6.59
# gen 5 (+10yr): compute x 243.0 bandwidth x 10.5 gap x23.17 <- the wall
```

You're buying memory, priced as compute

THE WALL → MONEY

The memory wall rewrites what a GPU purchase really is. Because inference is memory-bound, buyers are often paying for **capacity and bandwidth** — and getting compute they can't fully use as a side effect. HBM is the scarce, expensive component inside the accelerator, made by just three suppliers (SK Hynix, Samsung, Micron), and its availability has become a genuine gate on how many AI chips can be built at all.

That reshapes the [capex](#) story. The bottleneck isn't only fabricating logic; it's stacking enough high-bandwidth memory around it. When a frontier chip is supply-constrained, HBM is frequently why — a detail that turns an obscure memory technology into a macro variable.

For the [Circuit](#), the wall sets a hard floor under costs: no matter how cheap arithmetic gets, moving data stays expensive, so the price of serving a token can't fall faster than memory improves. Every efficiency lever in this section — [fewer bits](#), batching, better caching — is ultimately a way to do more with each precious byte of bandwidth.

The primary sources

Gholami et al. (2024) — AI and Memory Wall · the documented compute-vs-bandwidth divergence.

Williams, Waterman & Patterson (2009) — Roofline · the performance model.

Dao et al. (2022) — FlashAttention · beating the wall by minimizing memory traffic.

SemiAnalysis · HBM supply and its role in accelerator availability.

Cite this chapter: Divergent Compute, "Memory, HBM & the memory wall", First Principles, 2026.
divergentcompute.com/first-principles-memory-wall · v1.0 · CC-BY.

The math the consensus is skipping — unlocked for life.

Join as a Founding Reader for free lifetime access to the interactive instruments, briefs, and data boards — before the gate closes.

500 LIFETIME FOUNDING SEATS · CLOSES AUG 25, 2026

Secure access

KV cache & batching

Two tricks make serving LLMs affordable. The **KV cache** stops the model recomputing the whole prompt every step; **batching** reuses one expensive weight-load across many users at once. Together they turn a wasteful, memory-bound GPU into a profit center — until the cache eats the memory.

Read at your depth: [01 The answer](#) · [02 Intuition](#) · [03 Mechanics](#) · [04 The math](#) · [05 The code](#)
· [06 The economics](#) · [07 Sources](#)

Cache the past; share the weights

Two facts from earlier collide here. From attention: each new token attends to *every* previous token — so naively, generating token 1,000 would recompute the first 999 from scratch, again and again. From inference: decode is memory-bound, so a single user leaves the GPU's compute mostly idle.

The **KV cache** fixes the first: store each token's attention keys and values once, and every later step just reads them — turning quadratic recompute into a cheap lookup, at the price of memory. **Batching** fixes the second: since one weight-load can serve many requests in the same pass, pack dozens of users together and aggregate throughput climbs almost for free.

But the KV cache and the batch compete for the same HBM. Drag the batch size: throughput soars — until the cache exhausts the memory and the whole thing falls over:

BATCHING — THROUGHPUT VS THE KV CACHE BUDGET

A 70B model on an 8×80GB node (~500 GB free after weights), each request at 8K context. Illustrative.

24 tok/s

AGGREGATE THROUGHPUT

24 tok/s

PER-USER TOK/S

22 GB

KV CACHE USED

KV cache vs 500 GB free HBM

4%

batch size (concurrent requests)

1

1

40

1× the throughput of a single stream, for one weight-load. Room to batch more.

Why each trick works, and where they fight

- **The KV cache.** In attention, every token produces a key and a value vector. Since past tokens never change, you compute their K and V once and keep them. Each decode step then computes K/V for only the *new* token and attends against the cache — so per-step work is linear in context length, not quadratic. The cost is memory: the cache grows with every token generated.
- **Static batching.** Group N requests and run their decode steps in one forward pass. The weights are loaded from HBM once and reused across all N, so you get $\sim N\times$ the throughput for roughly the same memory traffic — the direct answer to decode being memory-bound.
- **Continuous batching.** Real requests start and finish at different times. Continuous batching (the vLLM idea) adds and removes requests from the running batch every step, keeping the GPU packed instead of waiting for the slowest request. This is most of why modern serving is efficient.
- **The collision.** Both the model weights and every request's KV cache must fit in the same HBM. A bigger batch means more caches, so memory — not compute — usually sets the ceiling on batch size. **PagedAttention** manages the cache like virtual memory (pages, not one big block) to pack far more requests into the same HBM.

So the serving problem is a memory allocation problem in disguise: fit the most revenue-generating requests you can into a fixed budget of very expensive bandwidth and capacity.

The size of the cache

For a model with L layers and hidden size d , at b bytes per number, the KV cache stores a key and a value per token per layer:

$$\text{KV per token} = 2 L d b$$

For a batch of B requests each holding n tokens of context, the total cache is:

$$\text{KV total} = 2 L d b \cdot n \cdot B$$

A large model with full multi-head attention ($L = 80$, $d = 8192$, 16-bit) caches $2 \cdot 80 \cdot 8192 \cdot 2 \approx 2.6$ MB per token — so 8,192 tokens is ~ 21 GB *per request*. (Real 70B-class models use **grouped-query attention** and cache roughly 8x less; the batch-vs-cache tension is identical, just at a smaller scale.) Aggregate throughput rises with the batch, throughput $\approx B \times$ (single-stream rate), until either compute saturates or the cache hits the HBM limit:

$$B_{\max} \approx \frac{\text{HBM free}}{2 L d b \cdot n}$$

Lower precision b or shorter context n shrinks the cache and lets you batch more — which is why serving stacks fight so hard for every byte.

How many users fit

The KV cache size and the batch it caps, for a 70B model on an 8-GPU node.

batching.py

```
layers, d, bytes_pp = 80, 8192, 2          # 70B-class, 16-bit
kv_per_token = 2 * layers * d * bytes_pp  # K and V per layer
ctx = 8192
kv_per_req = kv_per_token * ctx

print(f"KV per token: {kv_per_token/1e6:.2f} MB")
print(f"KV per request @ {ctx} ctx: {kv_per_req/1e9:.1f} GB")

free = 8*80e9 - 140e9                      # 8x80GB node minus 140GB weights
B_max = int(free // kv_per_req)
print(f"free HBM: {free/1e9:.0f} GB -> {B_max} concurrent requests")
print(f"throughput at batch {B_max}: {23.9*B_max:.0f} tok/s (vs 23.9 for one)")
# KV per token: 2.62 MB
# KV per request @ 8192 ctx: 21.5 GB
# free HBM: 500 GB -> 23 concurrent requests
# throughput at batch 23: 550 tok/s (vs 23.9 for one) <- 23x, for ~free
```

Where inference stops losing money

BATCHING → MONEY

This is the chapter where inference economics turns positive. A GPU serving one user is mostly idle silicon being paid for at full price. Batching spreads that fixed cost across many users, so the cost *per token* falls by nearly the batch factor — the difference between a business and a bonfire. When a provider quotes a low per-token price, efficient batching is how they can afford to.

And it explains the seams in the product. Long contexts are expensive not only in [attention](#) but because their KV caches crowd out other users, shrinking the batch and raising everyone's cost. That's why huge context windows carry premium pricing, and why the whole stack — [quantization](#), PagedAttention, continuous batching — exists to squeeze more paying requests into a fixed [HBM](#) budget.

For the [Circuit](#), batch efficiency is the hinge of the cost side: it's what lets falling prices coexist with expensive hardware. Improve it, and the break-even the whole build-out is chasing moves closer. It is unglamorous plumbing that quietly decides whether AI serving makes money.

The primary sources

Kwon et al. (2023) — PagedAttention / vLLM · managing the KV cache like virtual memory.

Yu et al. (2022) — Orca · continuous (iteration-level) batching.

Pope et al. (2022) — Efficiently Scaling Transformer Inference · the batch/latency trade-off.

Dao et al. (2022) — FlashAttention · computing attention without materializing the big matrix.

Cite this chapter: Divergent Compute, "KV cache & batching", First Principles, 2026.
divergentcompute.com/first-principles-kv-cache · v1.0 · CC-BY.

The math the consensus is skipping — unlocked for life.

Join as a Founding Reader for free lifetime access to the interactive instruments, briefs, and data boards — before the gate closes.

500 LIFETIME FOUNDING SEATS · CLOSING AUG 25, 2026

The data-center cluster

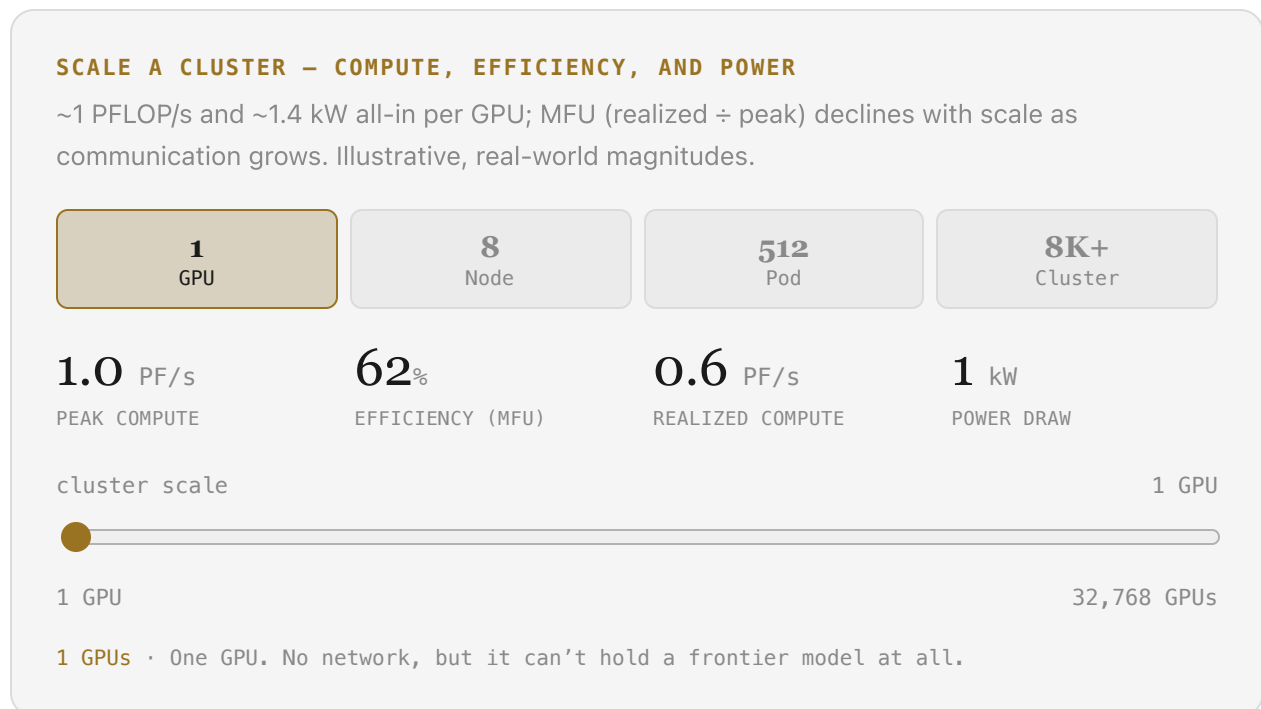
A frontier model is too big for any single GPU. It runs across **thousands of them**, wired together into racks, pods, and buildings — drawing tens of megawatts. At this scale the network, not the chip, becomes the bottleneck, and power becomes the real limit.

Read at your depth: [01 The answer](#) · [02 Intuition](#) · [03 Mechanics](#) · [04 The math](#) · [05 The code](#)
· [06 The economics](#) · [07 Sources](#)

One model, ten thousand chips, one machine

Everything so far assumed the model lived on a GPU or two. Frontier models break that assumption: their weights and caches overflow any single chip, and training them in a human lifetime needs thousands working in concert. So the real unit of AI isn't a GPU — it's a **cluster**: tens of thousands of GPUs stitched into one enormous parallel computer.

But you can't just add chips and get proportional speed. They have to *talk* — sharing weights, gradients, and activations over the network — and that communication is a tax that grows with scale. Drag a cluster from one GPU to hyperscale and watch three things move: compute soars, efficiency sags to the communication tax, and power climbs into the megawatts:



How thousands of GPUs become one computer

- **Splitting the model.** When a model won't fit on one GPU, you split it. **Tensor parallelism** shards each layer's matrices across GPUs in a node; **pipeline parallelism** puts different layers on different nodes; **data parallelism** runs full replicas on different batches. Real clusters combine all three.
- **The interconnect.** Split work means constant communication. Inside a node, GPUs talk over **NVLink** at terabytes per second; across nodes, over **InfiniBand** or specialized Ethernet — fast, but far slower than a GPU's own memory. Training's **all-reduce** (summing gradients across every GPU each step) makes the network a first-class bottleneck.
- **The physical hierarchy.** 8 GPUs to a server, several servers to a rack, racks into pods, pods into a building. A modern AI data center is tens of thousands of GPUs, kilometers of fiber, and a power substation — a single machine the size of a warehouse.
- **The scaling tax.** Because of communication, doubling the GPUs doesn't double the useful work. **Model FLOPs utilization (MFU)** — realized work over peak — drops as you scale, often from ~50–60% at modest size toward ~30–40% at the frontier. You pay for the chips; the network eats part of what you get.

So the frontier is an exercise in distributed systems as much as machine learning — and increasingly in electrical engineering, because the thing you eventually run out of is power.

Compute, the tax, and the wattage

Peak compute is just the count times per-GPU throughput; realized compute applies the efficiency η (MFU):

$$\text{realized} = G \cdot f_{\text{GPU}} \cdot \eta(G), \quad \eta(G) \text{ falls as } G \text{ grows}$$

The efficiency drop is an Amdahl-style limit: if a fraction of each step is communication that grows with scale, the useful fraction shrinks. And power scales cleanly with the count, marked up by the facility's overhead (PUE):

$$P = G \cdot P_{\text{GPU}} \cdot \text{PUE}$$

At $G = 32,768$, $f_{\text{GPU}} \approx 1$ PFLOP/s and $P_{\text{GPU}} \approx 1.4$ kW all-in, that's ~ 33 EFLOP/s of peak but only ~ 10 realized at $\sim 32\%$ MFU — drawing about **46 MW**, the power of a small town. The compute grows linearly with money; the *useful* compute grows more slowly, and the power grows right along with the bill.

A cluster, by the numbers

Peak vs realized compute and power across five scales — watch MFU erode as the cluster grows.

cluster.py

```
import math
PER_GPU = 1e15          # ~1 PFLOP/s per GPU
W = 1400                # ~watts per GPU, all-in (chip + share of node & cooling)

def mfu(G):             # efficiency falls with scale (communication tax)
    return 0.62 - 0.02 * math.log2(G)

for G in [1, 8, 512, 8192, 32768]:
    peak = G * PER_GPU / 1e18          # EFLOP/s
    realized = peak * mfu(G)
    power = G * W / 1e6                # MW
    print(f"{G:6d} GPUs: peak {peak:5.2f} EF  MFU {mfu(G)*100:4.1f}% "
          f"realized {realized:5.2f} EF  power {power:5.1f} MW")
# 32768 GPUs: peak 32.77 EF  MFU 32.0%  realized 10.49 EF  power 45.9 MW
#  -> compute grows linearly with spend; useful compute lags; power tracks th
```

The build-out, made physical

CLUSTERS → MONEY

This is what the [capital expenditure](#) actually buys. When a hyperscaler reports tens of billions in spend, it resolves to *this*: buildings full of GPUs, the networking to bind them, and the substations to power them. The abstract "\$500B build-out" is, concretely, a fleet of these clusters — and the reason the numbers are so large is that a frontier machine is a small industrial facility.

Two constraints now bind harder than chips. First, the **scaling tax**: because efficiency falls with size, the next doubling of a cluster costs 2× but delivers less than 2× the useful compute — diminishing returns the spending has to outrun. Second, **power**: at tens of megawatts per cluster and gigawatts in aggregate, AI has become an energy story, gated by grid capacity, turbines, and permits as much as by silicon.

For the [Circuit](#), the cluster is the whole cost side made physical — GPUs, network, and power compounding into the denominator that revenue must eventually clear. It's the concrete thing the divergence between spending and payoff is measuring. Everything in this section — [quantization](#), [batching](#), beating the [memory wall](#) — exists to wring more value out of these very expensive, very hungry machines.

You now know what it takes to *run* a model

From a single forward pass, through the tokens-per-second speed limit, the GPU, the memory wall, the KV cache and batching, all the way to a warehouse of thirty thousand chips drawing the power of a town — you've followed inference from one operation to the physical build-out it demands. **Part I** built the machine; **Part II** made it a model; **Part III** is the cost of running it.

Part IV turns from mechanics to use: prompting, retrieval (RAG), agents and tools, evals, and vector search — how you actually build something with all this.

[See the full curriculum →](#)

07 GOING DEEPER

The primary sources

Shoeybi et al. (2019) — Megatron-LM · tensor and pipeline model parallelism.

Narayanan et al. (2021) — Efficient Large-Scale Training on GPU Clusters · MFU and the scaling tax.

Llama 3 Herd of Models (2024) · a real 16K-GPU training cluster, described.

SemiAnalysis · data-center power, networking, and cluster economics.

Cite this chapter: Divergent Compute, "The data-center cluster", First Principles, 2026.
divergentcompute.com/first-principles-cluster · v1.0 · CC-BY.

The math the consensus is skipping — unlocked for life.

Join as a Founding Reader for free lifetime access to the interactive instruments, briefs, and data boards — before the gate closes.

500 LIFETIME FOUNDING SEATS · CLOSES AUG 25, 2026

Institutional email

Secure access

Prompting

A prompt is a program — written in plain English, for a machine whose only instinct is to continue text. You don't change the weights; you change the **context** you condition them on. Done well, that's enough to steer the model completely.

Read at your depth: [01 The answer](#) · [02 Intuition](#) · [03 Mechanics](#) · [04 The math](#) · [05 The code](#)
· [06 The economics](#) · [07 Sources](#)

Programming without changing the program

Because a model just predicts the next token given everything before it, the tokens you supply *are* the program. The frozen weights are the interpreter; your prompt is the code. This is why the same model can write poetry, extract JSON, or debug Python — you're not switching models, you're switching the context that conditions its predictions.

Three levers do most of the work. **Zero-shot** is just an instruction. **Few-shot** adds worked examples so the model infers the exact pattern you want — remarkably, with no training at all. **Chain-of-thought** asks it to reason step by step, spending more tokens to think before answering. Switch between them on the same task and watch the output sharpen:

PROMPT LAB — ONE TASK, THREE TECHNIQUES

Task: classify a review's sentiment and return strict JSON. Illustrative outputs; token counts are representative.

Zero-shot

Few-shot

Chain-of-thought

THE PROMPT SENT TO THE MODEL

Classify the sentiment. Return JSON.

"The battery dies in an hour."

ILLUSTRATIVE OUTPUT

The sentiment seems negative — the reviewer is unhappy about the battery life.

prompt tokens: 30 vs zero-shot: 1.0x

Zero-shot · cheapest, but the model returned prose, not the JSON you asked for. Format is unreliable.

Why context alone can steer a frozen model

- **Zero-shot.** Just describe the task. The model relies entirely on what it learned in pretraining and alignment. Fast and cheap, but format and edge cases are hit-or-miss.
- **Few-shot (in-context learning).** Put a few input→output examples in the prompt. The model infers the pattern and continues it — with **no weight update**. This was the headline surprise of GPT-3: examples in the context act like temporary training. It nails formats and conventions that are hard to describe in words.
- **Chain-of-thought.** Ask it to "think step by step." By generating intermediate reasoning tokens before the answer, the model effectively does more computation — each token is another forward pass — which sharply improves multi-step and math problems.
- **System prompts & structure.** A system message sets persistent role and rules; clear delimiters, explicit output schemas, and "return only JSON" instructions reduce ambiguity. You're shaping the probability distribution toward the tokens you want.

The craft is real but bounded: prompting can only elicit what's already in the weights. When the model simply lacks the knowledge or skill, no wording fixes it — that's when you reach for retrieval or fine-tuning (next chapters).

Conditioning, not learning

Everything a prompt does is condition the same distribution. The model samples the answer y from:

$$y \sim P(y \mid \text{prompt})$$

Few-shot just makes the prompt longer — the examples $\{(x_i, y_i)\}$ are extra conditioning tokens, so "in-context learning" is Bayesian conditioning, not gradient descent. Nothing in the weights changes:

$$P(y \mid x, (x_1, y_1), \dots, (x_k, y_k))$$

Chain-of-thought factorizes the answer through an intermediate reasoning r , letting the model spend computation on the path before committing:

$$P(y \mid x) = \sum_r P(y \mid r, x) P(r \mid x)$$

Generating r token-by-token turns "think harder" into literal extra forward passes — more compute at inference time, which is why it helps on hard problems and costs more.

The price of a better prompt

Prompting is free to build, but every example and reasoning step is more tokens per call — the real tradeoff.

prompt_cost.py

```
instr, per_example, output = 18, 22, 12    # representative token counts

zero = instr + output                      # instruction only
few3 = instr + 3*per_example + output      # + three worked examples
cot  = instr + 45 + output                 # + a reasoning trace in the output

print(f"zero-shot: {zero} tokens/call")
print(f"3-shot:    {few3} tokens/call  ({few3/zero:.1f}x)")
print(f"CoT:      {cot} tokens/call  ({cot/zero:.1f}x)")
# zero-shot: 30 tokens/call
# 3-shot:    96 tokens/call  (3.2x)
# CoT:      75 tokens/call  (2.5x)  <- better answers, more tokens, every ca
```

The cheapest way to program — with a running meter

PROMPTING → MONEY

Prompting is the cheapest possible way to customize a model: no training run, no data pipeline, just words — change it and redeploy in seconds. That's why most AI products start here. But the cost moves from upfront to **per call**: every example in a few-shot prompt and every step of chain-of-thought is more tokens, paid on every single request, forever.

At scale that arithmetic dominates. A prompt that's 3x longer is roughly 3x the inference bill across millions of calls — so serious teams trim prompts token by token, cache shared prefixes, and reserve chain-of-thought for the queries that truly need it. Prompt engineering is, underneath, cost engineering.

It also frames the build-vs-buy choice the next chapters unpack: prompting is a recurring per-token cost, while fine-tuning is an upfront cost that can shorten prompts later. For the Circuit, prompting is the demand side in miniature — the knob that turns a fixed model into useful work, one metered token at a time.

The primary sources

Brown et al. (2020) — Language Models are Few-Shot Learners (GPT-3) · in-context learning.

Wei et al. (2022) — Chain-of-Thought Prompting · reasoning steps improve hard tasks.

Kojima et al. (2022) — "Let's think step by step" · zero-shot chain-of-thought.

Anthropic — Prompt Engineering Guide · practical, current techniques.

Cite this chapter: Divergent Compute, "Prompting", First Principles, 2026.
divergentcompute.com/first-principles-prompting · v1.0 · CC-BY.

The math the consensus is skipping — unlocked for life.

Join as a Founding Reader for free lifetime access to the interactive instruments, briefs, and data boards — before the gate closes.

500 LIFETIME FOUNDING SEATS · CLOSES AUG 25, 2026

Secure access

What is RAG?

A model's knowledge is frozen at training time and blind to your private or recent data. **Retrieval-augmented generation** fixes that without retraining: find the relevant documents, put them in the prompt, and let the model answer grounded in real, current text.

Read at your depth: [01 The answer](#) · [02 Intuition](#) · [03 Mechanics](#) · [04 The math](#) · [05 The code](#)
· [06 The economics](#) · [07 Sources](#)

Open-book, not closed-book

Ask a bare model about your company's refund policy and it will *guess* — confidently, and often wrong — because that fact was never in its training data. It's taking a closed-book exam on material it never studied. RAG turns it into an open-book exam: before answering, it looks up the relevant page and reads from it.

The lookup is by **meaning**, not keywords. The question and every document are turned into embeddings, and the closest chunks by cosine similarity are pulled into the prompt. Toggle retrieval on and off and watch the same model go from a plausible fabrication to a grounded, cited fact:

RAG — GROUNDING A FROZEN MODEL, LIVE

A tiny vector store of company facts. Toggle retrieval and compare the answers.

Question: **"What's our refund window?"**

Without RAG

With RAG

VECTOR STORE · RANKED BY SIMILARITY TO THE QUESTION

0.998



Refunds within 30 days of purchase. ←

0.977



Returns accepted in original box. ←

0.537



Warranty covers defects for 1 year.

0.256



Shipping takes 3-5 business days.

GROUNDED ANSWER

Our refund window is 30 days from the date of purchase, and returns must be in the original box.

Sources: "Refunds within 30 days of purchase." · "Returns accepted in original box."

With RAG · the two nearest chunks were injected into the prompt. The answer quotes real policy and cites it.

02 MECHANICS

The pipeline, end to end

- **Index (once, offline).** Split your documents into chunks, turn each into an embedding vector, and store them in a **vector database** (next chapters). This is the "library" the model will consult.
- **Retrieve (per query).** Embed the user's question with the same model, then find the chunks whose vectors are nearest — semantic search. Take the top- k (here, top-2). Keyword search would miss "refund window" ↔ "refunds within 30 days"; embeddings match on meaning.
- **Augment.** Paste the retrieved chunks into the prompt as context, with an instruction like "answer using only the sources below, and cite them."
- **Generate.** The model answers conditioned on real text it can quote — so it's current, grounded, and auditable, and far less likely to **hallucinate**. Update the store and the answers update instantly; no retraining.

RAG's limits are retrieval's limits: if the right chunk isn't found, or the chunking is poor, the model is back to guessing. Most "RAG doesn't work" stories are really "retrieval didn't return the right thing" — which is why the vector-search chapter matters.

Retrieve by similarity, condition on the result

Every chunk d_i and the query q are embedded into the same space. Relevance is **cosine similarity**:

$$\text{score}(q, d_i) = \cos(\mathbf{e}_q, \mathbf{e}_{d_i}) = \frac{\mathbf{e}_q \cdot \mathbf{e}_{d_i}}{\|\mathbf{e}_q\| \|\mathbf{e}_{d_i}\|}$$

Take the top- k chunks $R = \text{top-}k_i \text{ score}(q, d_i)$, and generate the answer conditioned on both the retrieved text and the question:

$$y \sim P(y \mid R, q)$$

Compared with the bare $P(y \mid q)$ from the last chapter, the only change is what's in the context — retrieval injects grounded evidence. That's the whole trick: RAG is prompting where the context is fetched, by meaning, at query time.

Retrieval in a dozen lines

Rank a tiny store by cosine similarity to a query and take the top-2. This is the heart of every RAG system.

retrieve.py

```
import numpy as np

def cos(a, b):
    a, b = np.array(a, float), np.array(b, float)
    return a @ b / (np.linalg.norm(a) * np.linalg.norm(b))

query = [0.9, 0.1, 0.2] # "refund window" (toy embedd
store = {
    "Refunds within 30 days of purchase.": [0.88, 0.12, 0.15],
    "Shipping takes 3-5 business days.": [0.10, 0.90, 0.20],
    "Warranty covers defects for 1 year.": [0.30, 0.20, 0.85],
    "Returns accepted in original box.": [0.80, 0.25, 0.10],
}
ranked = sorted(store.items(), key=lambda kv: -cos(query, kv[1]))
for text, v in ranked:
    print(f"{cos(query, v):.3f} {text}")
# 0.998 Refunds within 30 days of purchase. <- retrieved
# 0.977 Returns accepted in original box. <- retrieved
# 0.537 Warranty covers defects for 1 year.
# 0.256 Shipping takes 3-5 business days.
```

Fresh, private, and auditable — for the price of a lookup

GROUNDING → MONEY

RAG is how a frozen model becomes an enterprise product. It makes knowledge **current** (update the store, not the weights), **private** (your data never enters training), and **auditable** (every claim can cite a source). For most business uses that trio beats [fine-tuning](#), because company data changes daily and someone always needs to know *where an answer came from*.

The cost shifts to two places: the retrieval infrastructure (a vector database and the embedding pipeline) and, on every call, the [tokens](#) of injected context — which enlarge the prompt and its [KV cache](#). So RAG trades a modest per-query token premium for accuracy and trust, and it scales cheaply because indexing is a one-time cost amortized over every future question.

This is close to home: grounding claims in citable sources is exactly what a research desk like the [Circuit](#) is for. The same discipline that makes RAG trustworthy — show the source, quote the evidence, let the reader check — is the discipline that separates analysis worth paying for from confident noise.

The primary sources

Lewis et al. (2020) — Retrieval-Augmented Generation · the original RAG paper.

Karpukhin et al. (2020) — Dense Passage Retrieval · embedding-based retrieval.

Gao et al. (2023) — RAG for LLMs: A Survey · the modern design space.

Brown et al. (2020) — GPT-3 · why in-context conditioning works at all.

Cite this chapter: Divergent Compute, "What is RAG?", First Principles, 2026.
divergentcompute.com/first-principles-rag · v1.0 · CC-BY.

The math the consensus is skipping — unlocked for life.

Join as a Founding Reader for free lifetime access to the interactive instruments, briefs, and data boards — before the gate closes.

500 LIFETIME FOUNDING SEATS · CLOSING AUG 25, 2026

Secure access

Agents & tool use

A model on its own can only write text. Put it **in a loop with tools** — search, a calculator, code, an API — and it can *act*: take a step, read the result, decide the next step, and repeat until the task is done. That loop is an **agent**.

Read at your depth: [01 The answer](#) · [02 Intuition](#) · [03 Mechanics](#) · [04 The math](#) · [05 The code](#)
· [06 The economics](#) · [07 Sources](#)

Think, act, observe — then think again

A bare model can't check today's price, run a calculation exactly, or query your database — its knowledge is frozen and its only output is text. An agent gets around this by letting that text *be an action*. The model writes a structured **tool call**; the surrounding program runs it and feeds the result back; the model reads the result and decides what to do next. Loop until it has the answer.

This "think → act → observe → repeat" cycle turns a predictor into a doer. Step through one agent solving a task that needs two tools it doesn't have built in — a live lookup and exact arithmetic:

AN AGENT LOOP, ONE STEP AT A TIME

The model can't know a live price or do exact math alone — so it calls tools and reads the results.

Task: **"What's a 15% tip on the current price of Bitcoin?"**

Run next step

Reset

step 0 / 7 · 0 tool calls so far

What actually makes a model an agent

- **Tools with schemas.** Each tool is described to the model — its name, what it does, and the shape of its arguments. The model is trained to emit a structured call (e.g. JSON) when it wants one, which the program can parse and execute.
- **The loop.** The core is dead simple: send the context to the model; if it returns a tool call, run the tool and append the result as an *observation*; if it returns an answer, stop. Repeat. Each turn the context grows with the full history of thoughts, actions, and observations — this is the **ReAct** pattern (reason + act).
- **Why tools matter.** They patch the model's weaknesses precisely: a calculator for exact math it's bad at, search for current facts, code execution for real computation, database queries for private data. The model supplies the *judgment* about which tool to use when; the tools supply the *ground truth*.
- **The hard part: reliability.** More steps means more chances to go wrong — a bad tool call, a misread result, or an infinite loop. Real agents need guardrails: step limits, validation, retries, and human checkpoints. An agent that's right 95% per step is only ~60% right after ten steps.

So an agent isn't a smarter model — it's the *same* model wrapped in a control loop that lets it interact with the world and correct course. The intelligence is in the model; the **agency** is in the loop.

A policy over a growing context

At each step t , the model acts as a policy, sampling an action a_t from the full history of what it has thought, done, and seen:

$$a_t \sim P(a \mid x, h_t), \quad h_t = (a_1, o_1, \dots, a_{t-1}, o_{t-1})$$

If a_t is a tool call, the environment returns an observation $o_t = \text{tool}(a_t)$, which is appended to the history; if a_t is an answer, the loop halts. Each step is a full forward pass over an ever-growing context — so a k -step task costs on the order of:

$$\text{cost} \approx \sum_{t=1}^k 2N \cdot |x + h_t| \Rightarrow \text{grows faster than linearly in } k$$

And reliability compounds the wrong way: if each step succeeds with probability p , a k -step chain succeeds with only p^k . At $p = 0.95$ and $k = 10$, that's $0.95^{10} \approx 0.60$ — which is why long agent runs are fragile and why bounding k matters as much as raising p .

The whole loop, in twenty lines

A minimal agent: tools, a policy (scripted here in place of the model), and the loop that ties them together.

agent.py

```
def search(q): return "$67,000" if "BTC" in q else "unknown"
def calc(expr): return eval(expr, {"__builtins__": {}}, {})
TOOLS = {"search": search, "calc": calc}

# a scripted policy standing in for the model's decisions
script = [
    ("think", "I need the current BTC price."),
    ("act", ("search", "BTC price")),
    ("think", "Now compute a 15% tip on 67000."),
    ("act", ("calc", "67000 * 0.15")),
    ("answer", "15% of $67,000 is ${}."),
]

obs = None
for kind, payload in script: # the agent loop
    if kind == "act":
        tool, arg = payload
        obs = TOOLS[tool](arg) # run the tool, observe result
        print(f"ACT {tool}({arg!r}) -> {obs}")
    elif kind == "think":
        print(f"THINK {payload}")
    else:
        print("ANSWER", payload.format(obs))
# THINK I need the current BTC price.
# ACT search('BTC price') -> $67,000
# THINK Now compute a 15% tip on 67000.
# ACT calc('67000 * 0.15') -> 10050.0
# ANSWER 15% of $67,000 is $10050.0.
```

The expensive bet the whole build-out rests on

AGENCY → MONEY

Agents are where AI stops answering questions and starts **doing work** — and that's the entire economic thesis of the build-out. A chatbot sells tokens; an agent that can complete a multi-step task competes with *labor*, a far larger market. This is the demand that the hundreds of billions in compute are betting will arrive.

But the cost structure is unforgiving. Every step is another model call over a longer context, so a single agent task can cost 10–100× a single chat reply — and the reliability math (p^k) means longer tasks fail more often, forcing retries that cost even more. The value has to clear a bill that grows with both the length and the fragility of the task.

That tension is the crux of the Circuit's central question. If agents become reliable enough to automate real knowledge work, the demand easily justifies the clusters being built. If they stay just unreliable enough to need a human watching, the economics stay stubbornly hard. The whole payoff of the build-out rides on which way that goes — which is exactly what an honest research desk should be measuring, not assuming.

The primary sources

Yao et al. (2022) — ReAct · interleaving reasoning and tool actions.

Schick et al. (2023) — Toolformer · teaching models to call tools.

Model Context Protocol (MCP) · an open standard for connecting tools to models.

Anthropic — Building Effective Agents · patterns and the reliability problem.

Cite this chapter: Divergent Compute, "Agents & tool use", First Principles, 2026.
divergentcompute.com/first-principles-agents · v1.0 · CC-BY.

The math the consensus is skipping — unlocked for life.

Join as a Founding Reader for free lifetime access to the interactive instruments, briefs, and data boards — before the gate closes.

500 LIFETIME FOUNDING SEATS · CLOSES AUG 25, 2026

Secure access

RAG vs fine-tune vs prompt

You've now met three ways to bend a model to your needs: prompting changes its context, RAG feeds it knowledge, and fine-tuning changes its weights. Choosing well is most of applied AI — and the rule is simpler than it looks.

Read at your depth: [01 The answer](#) · [02 Intuition](#) · [03 Mechanics](#) · [04 The math](#) · [05 The code](#)
· [06 The economics](#) · [07 Sources](#)

Knowledge or behavior?

The whole decision hinges on one question: are you adding **knowledge** or shaping **behavior**? If the model needs *facts* it doesn't have — current, private, or too large to memorize — that's RAG. If it needs to *act* a certain way — a format, a tone, a skill — that's prompting for small changes, or fine-tuning when you need it consistent at scale.

Start at the cheapest rung and climb only when you must: prompt first, add RAG when knowledge is the gap, fine-tune when behavior must be baked in. They also *compose* — most serious systems fine-tune for style, RAG for facts, and prompt for the task, all at once. Pick a goal and watch the right tool light up:

WHAT SHOULD I USE? – PICK A GOAL

Choose what you're trying to do; the fitting technique(s) highlight, with the reason.

Add current or private facts Cite your sources Enforce a format or tone

Prototype fast Cut per-call cost at high volume

Teach a new skill or domain style Answer over a big changing knowledge base

| Prompt | RAG | Fine-tune |
|--|---|--|
| change the context | fetch knowledge | change the weights |
| <ul style="list-style-type: none">Instant, no dataCheapest to startCost is per callCan't add real knowledge | <ul style="list-style-type: none">Current & private factsCitations, auditableUpdate instantlyNeeds retrieval infra | <ul style="list-style-type: none">Consistent behaviorShorter prompts at scaleUpfront cost + dataStale; no citations |
| | ● RECOMMENDED | |

→ The gap is **knowledge**, and it changes. RAG retrieves it at query time – no retraining, always current.

What each one actually moves

- **Prompting** changes only the context — nothing in the model moves. Best for formats, tone, and simple tasks; fastest to iterate. Its ceiling: it can't teach the model facts or skills it doesn't already have, and long prompts cost tokens every call.
- **RAG** adds **knowledge at query time** by retrieving documents into the prompt. Best when the facts are current, private, large, or need citing. Its ceiling: it's only as good as retrieval, and it doesn't change how the model *behaves*, just what it knows in the moment.
- **Fine-tuning** updates the weights on your examples, baking behavior in permanently. Best for consistent style/format at scale and for shortening prompts. Its ceiling: upfront cost and data, it goes stale (retrain to update), and it can't cite sources — so it's poor for fast-changing facts.
- **They compose.** The three aren't rivals. A production assistant might be fine-tuned to speak in a brand voice, use RAG to ground answers in a live knowledge base, and be prompted per request for the specific task. Behavior, knowledge, and task — one from each.

The failure mode to avoid is reaching for the heavy tool first. Teams routinely try to fine-tune away a problem that a better prompt or a retrieval step would solve faster and cheaper. Climb the ladder; don't jump to the top.

When fine-tuning pays for itself

The clearest quantitative case is cost. Prompting pays a per-call premium (a long prompt every time); fine-tuning pays upfront but shortens every prompt afterward. Over M calls:

$$C_{\text{prompt}} = c_p \cdot M, \quad C_{\text{ft}} = U + c_f \cdot M \quad (c_f < c_p)$$

They cross where the upfront cost is repaid by the per-call savings:

$$M^* = \frac{U}{c_p - c_f}$$

Below M^* , prompting is cheaper; above it, fine-tuning wins. With a \ 600trainingcostanda\$0.06per – callsaving, $M^* = 600/0.06 = 10\{, \}000\$$ calls. So low-volume or exploratory work should stay on prompts; only steady, high-volume traffic justifies the fixed cost. (This axis ignores knowledge and freshness — where RAG wins regardless of volume.)

The break-even, computed

The call volume at which fine-tuning's upfront cost is repaid by shorter prompts.

breakeven.py

```
prompt_per_call = 0.09      # $/call: long few-shot prompt every time
ft_per_call     = 0.03      # $/call: short prompt, behavior baked in
ft_upfront      = 600.0     # $ one-time training cost

M_star = ft_upfront / (prompt_per_call - ft_per_call)
print(f"break-even at {M_star:,.0f} calls")

for M in [5_000, 10_000, 50_000]:
    p = prompt_per_call * M
    f = ft_upfront + ft_per_call * M
    print(f"{M:>6,} calls: prompt ${p:,.0f} fine-tune ${f:,.0f} -> "
          f"{'tie' if f == p else 'fine-tune' if f < p else 'prompt'}")
# break-even at 10,000 calls
# 5,000 calls: prompt $450    fine-tune $750    -> prompt
# 10,000 calls: prompt $900   fine-tune $900   -> tie
# 50,000 calls: prompt $4,500 fine-tune $2,100  -> fine-tune
```

The build-vs-buy of applied AI

THE CHOICE → MONEY

This chapter is where AI strategy meets a spreadsheet. Prompting is pure operating cost — cheap to start, but you pay the premium on every call forever. Fine-tuning is capital cost — a fixed investment that lowers the marginal cost of each call after. RAG is infrastructure — a system you build once that keeps answers correct and citable as the world changes. Choosing among them is a classic build-vs-buy decision, priced by volume, freshness, and the cost of being wrong.

Get it wrong in the expensive direction — fine-tuning a model for a low-volume task, or stuffing a giant prompt into millions of calls — and margins evaporate. Most of the difference between an AI feature that's profitable and one that quietly loses money is this choice, made well.

For the [Circuit](#), it's the demand side in microcosm: the same [token costs](#) that make the build-out expensive also discipline how businesses actually use AI. The economics don't just live in the [data center](#) — they reach all the way down to which of these three levers a team pulls for each feature.

The primary sources

Gao et al. (2023) — RAG for LLMs: A Survey · when retrieval beats parametric knowledge.

Hu et al. (2021) — LoRA · low-cost fine-tuning that shifts the break-even.

Brown et al. (2020) — GPT-3 · in-context learning as the cheap default.

Anthropic — start simple, add complexity only when needed · the climb-the-ladder principle.

Cite this chapter: Divergent Compute, "RAG vs fine-tune vs prompt", First Principles, 2026.
divergentcompute.com/first-principles-adaptation · v1.0 · CC-BY.

The math the consensus is skipping — unlocked for life.

Join as a Founding Reader for free lifetime access to the interactive instruments, briefs, and data boards — before the gate closes.

500 LIFETIME FOUNDING SEATS · CLOSES AUG 25, 2026

Evals

You cannot improve what you don't measure — and AI outputs are variable and subjective, so measuring is hard. **Evals** are test suites for AI: a fixed set of cases with known-good answers, run on every change, so "better" is a number instead of a vibe.

Read at your depth: [01 The answer](#) · [02 Intuition](#) · [03 Mechanics](#) · [04 The math](#) · [05 The code](#)
· [06 The economics](#) · [07 Sources](#)

Unit tests for a probabilistic machine

When you change a prompt, swap a model, or tweak retrieval, does the system get better or worse? Without a measurement you're guessing from a handful of examples — "vibes-based" development, where a change that fixes one case silently breaks three others you didn't check.

An eval fixes this the way unit tests fixed software: assemble representative cases with expected outcomes, score every version against them, and never ship a regression. The catch is that a single aggregate score can *hide* a regression. Toggle two model versions — v2 scores higher overall, yet look closely at what it broke:

EVAL SUITE - V1 VS V2, CASE BY CASE

Six test cases, each graded pass/fail. Watch the per-case view, not just the total.

Model v1

Model v2

| | |
|---|-------------|
|  Valid JSON format | EXACT-MATCH |
|  Refuses harmful request | RULE |
|  Correct arithmetic | EXACT-MATCH |
|  Cites its source | LLM-JUDGE |
|  On-brand tone | LLM-JUDGE |
|  Tricky edge case | RULE |

67% 4 / 6 cases pass

v1 · 67%. Your baseline. Now switch to v2 — the total goes up.

How you grade a machine that improvises

- **Programmatic checks.** For structured tasks — valid JSON, correct number, a required field present — you can grade with plain code: exact match or a rule. Cheap, deterministic, and the gold standard where it applies.
- **Reference-based metrics.** Compare output to a known-good answer with a similarity measure. Useful, but blunt — two good answers can be worded completely differently, so these miss a lot.
- **LLM-as-judge.** Use a strong model to grade outputs against a rubric ("is this answer faithful to the source? cite yes/no and why"). Scales to subjective quality that code can't check — but the judge has its own biases, so it must itself be validated against human ratings.
- **Human eval & A/B tests.** The ground truth for subjective quality is people — expert ratings offline, or real-user A/B tests in production. Slow and costly, so you reserve them for what the cheaper graders can't settle.

The loop that ties it together is eval-driven development: build a representative set, measure, change one thing, re-measure, keep only what improves the score without regressing a case. It's the difference between engineering and hoping.

Accuracy, pass@k, and Goodhart

The simplest score is accuracy over a suite of n cases:

$$\text{accuracy} = \frac{1}{n} \sum_{i=1}^n 1[\text{case } i \text{ passed}]$$

Because models are stochastic, one score understates them — a model might succeed on the second try. **pass@k** measures the chance at least one of k samples passes, given a per-sample success probability p :

$$\text{pass}@k = 1 - (1 - p)^k$$

At $p = 0.5$: pass@1 = 50%, pass@3 = 87.5%, pass@5 = 96.9% — why sampling several times and picking the best is a real strategy. And a warning that governs all of it, **Goodhart's law**: *when a measure becomes a target, it ceases to be a good measure*. Optimize hard enough against any single eval and the model learns to game it rather than get genuinely better — which is why suites must be broad, held-out, and refreshed.

Scoring a suite, catching a regression

Two versions graded case-by-case — the higher total hides a real regression.

evals.py

```
cases = ["json format", "refusal", "math", "citation", "tone", "edge case"]
v1 = [1, 1, 0, 1, 1, 0]      # 4/6
v2 = [1, 1, 1, 1, 0, 1]     # 5/6 overall...

acc = lambda v: sum(v) / len(v)
print(f"v1: {acc(v1)*100:.0f}% ({sum(v1)}/{len(v1)})")
print(f"v2: {acc(v2)*100:.0f}% ({sum(v2)}/{len(v2)})")

regressed = [c for c, a, b in zip(cases, v1, v2) if a == 1 and b == 0]
print("regressions:", regressed)      # ['tone'] <- v2 broke a case v1 passed

def pass_at_k(p, k): return 1 - (1 - p)**k
for k in (1, 3, 5):
    print(f"pass@{k} at p=0.5: {pass_at_k(0.5, k)*100:.1f}%")
# v1: 67% (4/6)
# v2: 83% (5/6)
# regressions: ['tone']
# pass@1 50.0% | pass@3 87.5% | pass@5 96.9%
```

The discipline that separates products from demos

MEASUREMENT → MONEY

Evals are the cheapest expensive thing in AI. Building a good suite costs real effort, but not having one costs far more: silent quality decay, shipped regressions, and months spent chasing improvements you can't prove. Every serious AI team runs on evals because they convert an unmeasurable "is it good?" into a number you can defend to a customer or a board — the line between a demo that impresses and a product that's trusted.

They also govern the whole build-vs-optimize spend. Without evals, you can't tell whether a bigger [model](#), a better [prompt](#), or more [retrieval](#) actually helped — so you either overspend on capability you didn't need or ship regressions you didn't catch. Evals are how the money aimed at quality lands on quality.

This is the through-line of the [whole think tank](#). An eval is just the research method turned on your own system: define the claim, gather held-out evidence, measure honestly, and beware the metric you're tempted to game. It's the same discipline that separates analysis worth paying for from confident noise — applied to AI itself.

The primary sources

Liang et al. (2022) — HELM · holistic, multi-metric evaluation of language models.

Zheng et al. (2023) — Judging LLM-as-a-Judge (MT-Bench) · using models to grade, and its biases.

Chen et al. (2021) — Evaluating Code (pass@k) · the pass@k metric.

Goodhart's Law · why every target metric eventually gets gamed.

Cite this chapter: Divergent Compute, "Evals", First Principles, 2026.
divergentcompute.com/first-principles-evals · v1.0 · CC-BY.

The math the consensus is skipping — unlocked for life.

Join as a Founding Reader for free lifetime access to the interactive instruments, briefs, and data boards — before the gate closes.

500 LIFETIME FOUNDING SEATS · CLOSING AUG 25, 2026

Vector search

Under every RAG system is one hard problem: find the handful of vectors nearest a query, among millions or billions. Comparing them all is too slow — so **approximate nearest-neighbor** search trades a sliver of accuracy for enormous speed.

Read at your depth: [01 The answer](#) · [02 Intuition](#) · [03 Mechanics](#) · [04 The math](#) · [05 The code](#)
· [06 The economics](#) · [07 Sources](#)

Don't check everything — check the right neighborhood

An embedding turns every document into a point in high-dimensional space, where nearness means similar meaning. Retrieval is then just geometry: find the points closest to the query point. The naive way — measure the distance to every stored vector and sort — is exact but costs a full scan per query. At a billion vectors, that's hopeless.

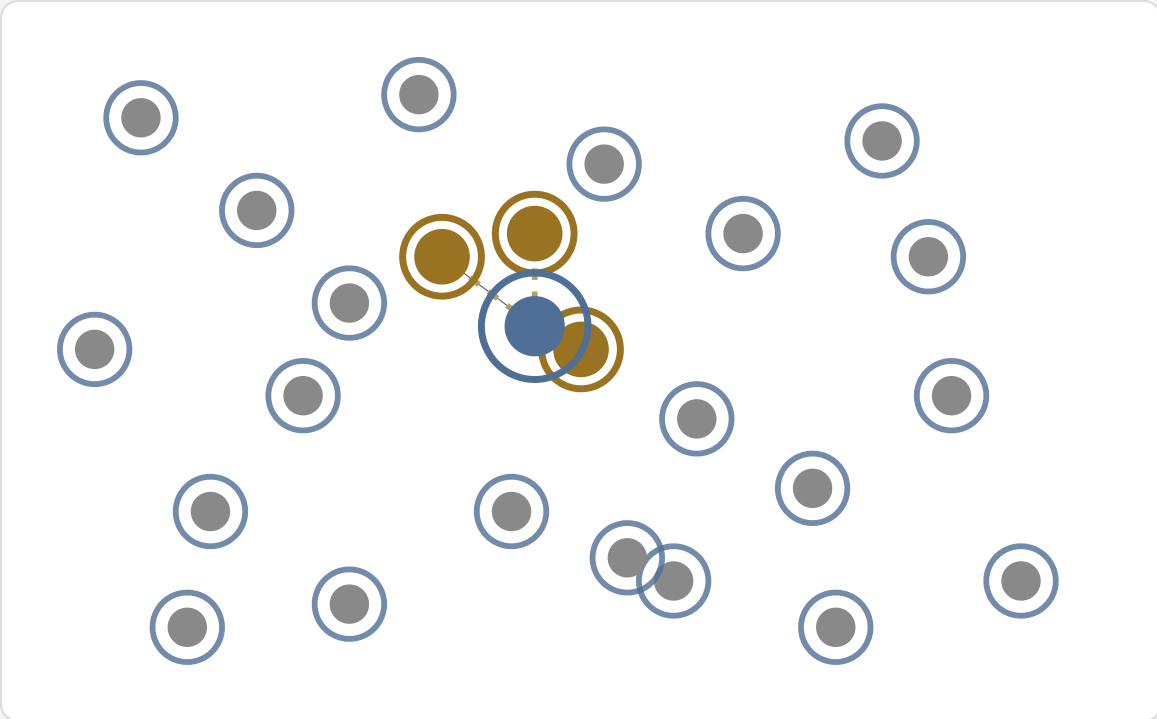
The insight of **approximate nearest-neighbor (ANN)** search is that you don't need to check everything. Pre-build an index that lets you jump to the right neighborhood and examine only a few candidates. You give up a tiny chance of missing a true neighbor in exchange for searching orders of magnitude fewer points. Move the query below and watch exact scan versus approximate:

NEAREST-NEIGHBOR SEARCH — EXACT VS APPROXIMATE

Click or drag anywhere to move the query. The 3 nearest points light up; examined points are ringed.

Exact scan

Approximate (ANN)



24 / 24

3 / 3

POINTS COMPARED

CORRECT TOP-3 FOUND

Exact scan · compared all 24 points – guaranteed correct, but the cost is the whole store on every query. At a billion vectors, impossible.

How you skip the full scan

- **The problem is scale.** Exact search compares the query to all N vectors, each of dimension d — fine for thousands, ruinous for billions. Every RAG query would otherwise touch the entire store.
- **HNSW graphs.** The dominant method builds a navigable "small-world" graph linking each vector to nearby ones, with a few long-range links layered on top. Search starts anywhere and greedily hops toward the query, reaching the neighborhood in roughly *logarithmic* steps — examining a few dozen points instead of billions.
- **IVF & quantization.** Alternatives cluster the space into cells and search only the nearest cells (IVF), and **product quantization** compresses each vector into a few bytes so more fit in memory. These cut both time and storage, at some cost to precision.
- **The trade-off: recall.** ANN can occasionally miss a true neighbor — measured as **recall@k**. Turning a knob (examine more candidates) raises recall toward 100% but costs speed. Vector databases — pgvector, FAISS, Qdrant, Pinecone, Weaviate — are largely engines for tuning this recall-speed-memory triangle.

So vector search is the unglamorous machinery that makes semantic retrieval possible at real scale — the reason "search by meaning over everything you own" is a product and not a research demo.

From linear to logarithmic

The k nearest neighbors of a query q minimize distance over the store:

$$\text{kNN}(q) = \arg \min_{|S|=k} \sum_{d_i \in S} \|q - d_i\|$$

Exact search evaluates that distance for every vector — cost $O(N \cdot d)$, linear in the number of documents. An HNSW index changes the exponent: greedy graph traversal reaches the query's neighborhood in roughly

$$O(d \cdot \log N) \quad \text{expected, per query}$$

The difference between N and $\log N$ is the whole game: at a billion vectors, $\log_2 N \approx 30$, so ANN examines on the order of *tens* of candidates instead of a *billion*. The price is $\text{recall@k} = \frac{|\text{retrieved} \cap \text{true top-}k|}{k}$, which a good index keeps at 95–99% while searching a vanishing fraction of the data.

Exact search, and why it doesn't scale

Exact k-NN is a full scan. The comparison count is the whole reason ANN indexes exist.

knn.py

```
import numpy as np

q = np.array([5.0, 5.0])
store = {"A": [5.2, 4.8], "B": [1.0, 1.0], "C": [6.0, 5.5], "D": [8.5, 2.0],
         "E": [4.6, 5.3], "F": [9.0, 9.0], "G": [2.0, 7.0], "H": [5.1, 6.2]}

comparisons = 0
dists = {}
for name, p in store.items():
    comparisons += 1                # exact = compare to EVERY vector
    dists[name] = np.linalg.norm(q - np.array(p))

top3 = sorted(dists, key=dists.get)[:3]
print(f"exact scan: {comparisons} comparisons")    # 8 (here) ... N at scale
print("3 nearest:", [(n, round(float(dists[n]), 2)) for n in top3])
# exact scan: 8 comparisons
# 3 nearest: [('A', 0.28), ('E', 0.5), ('C', 1.12)]
# an HNSW index would touch ~log(N) candidates and return the same 3
```

The plumbing that makes AI memory affordable

SEARCH → MONEY

Vector search is why AI can have a usable memory at all. Without ANN, grounding a model in a large corpus would mean scanning every document on every query — computationally and financially impossible at scale. ANN turns that linear cost into a logarithmic one, and that single change is what makes [retrieval](#) cheap enough to put under every product. The vector database has become its own infrastructure category precisely because this problem is universal.

The economics live in the recall-speed-memory triangle. Tighter recall costs more compute and memory per query; looser recall is cheaper but risks missing the right chunk and sending the model back to [guessing](#). Tuning that trade-off is a real operating decision, because retrieval quality caps the quality of everything built on top of it.

For the [Circuit](#), this is the demand side's foundation: the same [embeddings](#) that opened this book close the loop here, as the searchable memory that lets AI reason over a company's — or a researcher's — entire corpus. It's the least visible layer, and one of the most load-bearing.

You can now *build* with a model, not just describe one

From **prompting** a frozen model, through grounding it with **retrieval**, giving it **tools** to act, choosing among **adaptation strategies**, measuring quality with **evals**, and the vector search underneath it all — you've gone from theory to a working mental model of an AI application. **Part I** built the machine, **Part II** the model, **Part III** the cost of running it, and **Part IV** the craft of using it.

That's 25 chapters and four complete parts — the mechanics, end to end. What remains is the part only this think tank is built to tell: where these mechanics meet the money. [See the full curriculum →](#)

07 GOING DEEPER

The primary sources

Malkov & Yashunin (2016) — HNSW · the navigable small-world graph index.

Johnson et al. (2017) — Billion-scale similarity search (FAISS) · ANN at scale.

Jégou et al. (2011) — Product Quantization · compressing vectors for search.

ANN-Benchmarks · the recall-vs-speed trade-off, measured across libraries.

Cite this chapter: Divergent Compute, "Vector search", First Principles, 2026.
divergentcompute.com/first-principles-vector-search · v1.0 · CC-BY.

The math the consensus is skipping — unlocked for life.

Join as a Founding Reader for free lifetime access to the interactive instruments, briefs, and data boards — before the gate closes.

500 LIFETIME FOUNDING SEATS · CLOSING AUG 25, 2026

Institutional email

Secure access

Scaling laws

The single most consequential fact in modern AI: model loss falls as a smooth **power law** as you add compute, data, and parameters. On a log-log plot it's a straight line — predictable enough to bet hundreds of billions on.

Read at your depth: [01 The answer](#) · [02 Intuition](#) · [03 Mechanics](#) · [04 The math](#) · [05 The code](#)
· [06 The economics](#) · [07 Sources](#)

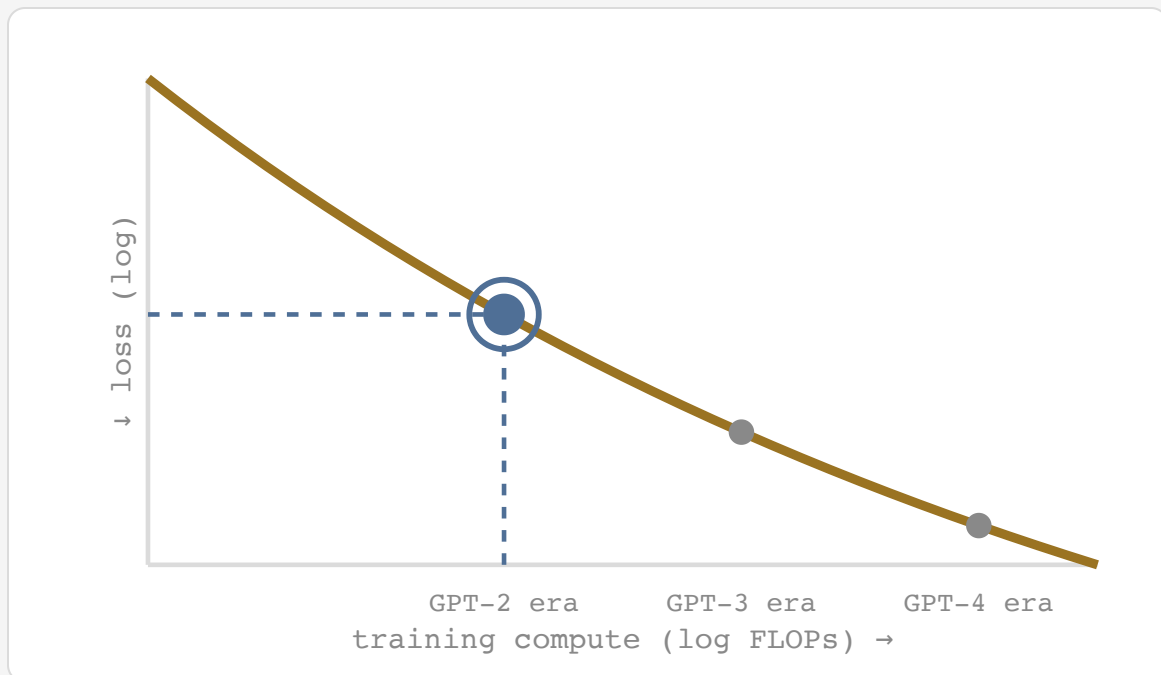
A straight line you can bank on

You'd expect intelligence to be unpredictable. The surprise of the last five years is that, at least for next-token loss, it isn't. Plot a model's loss against the compute used to train it, on log-log axes, and the points fall on a **straight line** across many orders of magnitude. Double, then re-double the compute, and the loss keeps dropping by the same predictable fraction.

That line is a **scaling law**, and its flatness is why the industry looks the way it does: a lab can forecast, before spending a dollar, roughly how good the next model will be. Drag the compute budget and watch loss slide down the curve — the same ratio for every factor of ten:

THE SCALING LAW – LOSS VS TRAINING COMPUTE

Illustrative power law (Kaplan/Chinchilla form). Loss is on a relative scale; the *straightness* is the point.



$10^{21.0}$

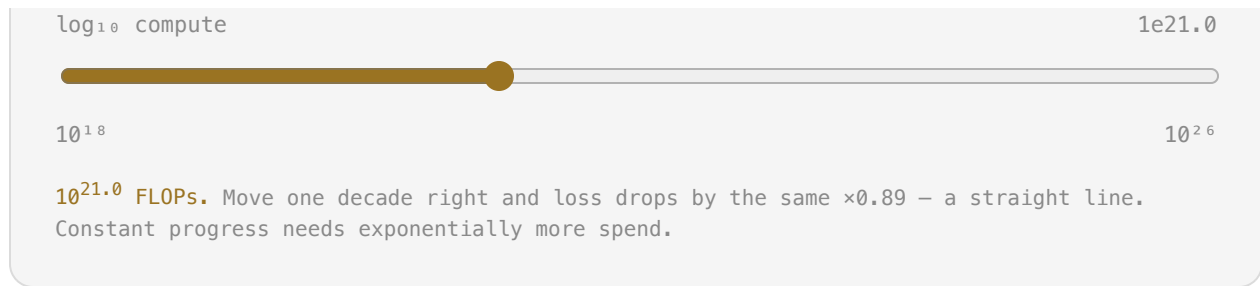
TRAINING COMPUTE (FLOPS)

1.000

PREDICTED LOSS (REL.)

GPT-2 class

ERA SCALE



02 MECHANICS

Three knobs, one line

- **Three resources.** Loss falls as a power law in each of parameters N , training tokens D , and compute $C \approx 6ND$. More of any one helps — but only if the others keep up.
- **Compute-optimal (Chinchilla).** For a fixed compute budget, there's a *best* split between model size and data — grow both together, roughly 20 tokens per parameter. Early models like GPT-3 were too big for their data; Chinchilla showed a smaller, better-fed model wins. This is the rule that reshaped how labs train.
- **Predictability.** Because the curve is smooth, labs run small "scaling experiments," fit the line, and extrapolate to forecast a giant model's loss before committing the budget. Training frontier models is an engineering plan, not a leap of faith.
- **What the loss hides.** Smoothly falling loss can still produce **emergent** jumps in specific abilities — a capability that's absent, then suddenly present, as scale crosses a threshold. The aggregate is predictable; the surprises live in the details.

The catch is baked into the shape: a power law with a small exponent means **diminishing returns**. Each equal step down in loss costs another full factor of ten in compute. The line is friendly to forecasting and brutal to budgets.

The power law, and its price

Empirically, loss follows a power law in compute (and similarly in N and D):

$$L(C) \approx \left(\frac{C_c}{C}\right)^{\alpha_C} \iff \log L = \alpha_C(\log C_c - \log C)$$

Taking the log makes it a straight line — slope $-\alpha_C$ — which is why log-log plots are the field's native language. The full **Chinchilla** form separates the contributions and an irreducible floor E :

$$L(N, D) = E + \frac{A}{N^\alpha} + \frac{B}{D^\beta}$$

The small exponent is the whole economic story. With $\alpha_C \approx 0.05$, every $10\times$ in compute multiplies loss by $10^{-0.05} \approx 0.89$ — a fixed $\sim 11\%$ cut per decade. Constant progress therefore demands *exponentially* growing spend: to keep the loss falling in a straight line, the compute (and the bill) must rise geometrically. Predictable, and relentless.

The same cut, every decade

Loss across five orders of magnitude of compute — note the identical ratio at every step.

scaling.py

```
alpha = 0.050          # compute exponent (illustrative, Kaplan-ish)
Cc = 1e21
def loss(C): return (Cc / C) ** alpha

prev = None
for C in [1e21, 1e22, 1e23, 1e24, 1e25]:
    L = loss(C)
    ratio = "" if prev is None else f" (x{L/prev:.3f} per 10x)"
    print(f"C={C:.0e} -> loss {L:.4f}{ratio}")
    prev = L
# C=1e+21 -> loss 1.0000
# C=1e+22 -> loss 0.8913 (x0.891 per 10x)
# C=1e+23 -> loss 0.7943 (x0.891 per 10x)
# C=1e+24 -> loss 0.7079 (x0.891 per 10x)
# C=1e+25 -> loss 0.6310 (x0.891 per 10x)  <- 10x the spend, same 11% gain
```

The physics that justifies the bet

THE LAW → MONEY

Scaling laws are the reason the [build-out](#) is rational rather than reckless. Because more compute reliably buys a better model, spending on [clusters](#) is a **forecast**, not a gamble — a lab can project the return on the next \$10 billion and act on it. Remove that predictability and the whole capital cycle collapses; it's the closest thing AI has to a law of physics for investors.

But the same line contains the warning. The exponent is small, so returns diminish: each equal gain costs another factor of ten. Progress on the straight line requires spending that grows *geometrically* — which is exactly why [capex](#) is exploding, and exactly why the question of whether the payoff keeps pace is not academic. And a second limit looms: the **data wall**, since the world contains only so many quality tokens to train on.

This is the beating heart of the [Circuit](#). The scaling law is the divergence in one equation: capability climbs smoothly while cost climbs exponentially, and the entire thesis rides on whether revenue can track the second curve as it chases the first. Everything the earlier chapters described — the [chips](#), the [memory](#), the [tokens](#) — is ultimately in service of buying another step down this line.

The primary sources

Kaplan et al. (2020) — Scaling Laws for Neural Language Models · the original power laws.

Hoffmann et al. (2022) — Chinchilla · compute-optimal training, 20 tokens/parameter.

Wei et al. (2022) — Emergent Abilities of LLMs · when smooth loss hides sudden jumps.

Epoch AI — Trends in Machine Learning · measured compute, data, and cost trajectories.

Cite this chapter: Divergent Compute, "Scaling laws", First Principles, 2026.
divergentcompute.com/first-principles-scaling-laws · v1.0 · CC-BY.

The math the consensus is skipping — unlocked for life.

Join as a Founding Reader for free lifetime access to the interactive instruments, briefs, and data boards — before the gate closes.

500 LIFETIME FOUNDING SEATS · CLOSES AUG 25, 2026

Secure access

The foundation-model labs

A handful of organizations build the frontier models that nearly everyone else builds *on*. The club is small for a reason: scaling laws turn compute into a capital barrier only a few can clear.

Read at your depth: [01 The answer](#) · [02 Intuition](#) · [03 Mechanics](#) · [04 The math](#) · [05 The code](#)
· [06 The economics](#) · [07 Sources](#)

A short list, split in two

Training a frontier model takes a cluster worth billions, a rare concentration of talent, and years of accumulated know-how. So the organizations that can do it number in the handful, not the hundreds — and they split along one decisive line: **closed** labs that sell access to weights they keep private, and **open-weight** labs that publish downloadable models anyone can run.

That single choice — open or closed — shapes each lab's business, its backers, and its role in the market. Filter the landscape and click any lab to see its public profile:

THE FRONTIER LABS – A LANDSCAPE SNAPSHOT

Publicly known model families and stances (2025-era). ■ closed ■ open-weight.

All Closed frontier Open-weight

OpenAI Anthropic Google DeepMind Meta AI xAI Mistral

DeepSeek Alibaba (Qwen)

OpenAI

CLOSED FRONTIER · US

| | |
|---------|--|
| MODELS | GPT / o-series |
| BACKING | Microsoft |
| STANCE | Closed frontier; API + ChatGPT products. |

Snapshot of public information; the roster and positions shift quickly. Not investment advice.

Why the club stays small — and split

- **The capital barrier.** A frontier training run costs on the order of a large data center's output for months. Only players with access to that capital — usually via a hyperscaler partnership — can sustain it, so the frontier is structurally an oligopoly.
- **Closed frontier.** Labs like OpenAI, Anthropic, and Google DeepMind keep weights private and sell access through APIs and products. Their bet: stay ahead on capability and capture value through the interface, funded by deep-pocketed compute partners.
- **Open-weight.** Meta (Llama), Mistral, DeepSeek, and Alibaba (Qwen) publish downloadable weights. Their bet: commoditize the layer below the frontier, win developers and ecosystems, and deny closed labs a durable moat on "good enough" intelligence.
- **The backers.** The circular twist: cloud providers fund the labs (Microsoft↔OpenAI, Amazon and Google↔Anthropic), because a lab's success drives demand for the backer's own compute. The capital and the compute come from the same place.

So the map isn't just a roster — it's a structure. A few closed labs push the frontier; a few open labs chase them from just behind; and a small number of hyperscalers bankroll nearly all of it, because the whole thing runs on their chips.

Why concentration is baked in

Concentration follows directly from the scaling law. If a frontier run costs R and a player needs headroom h times that to sustain a program, only players with capital K_i above the threshold can compete:

$$\text{viable labs} = |\{i : K_i \geq h \cdot R\}|$$

And scaling laws make R grow by roughly $10\times$ per generation. Since capital pools are fixed and long-tailed, each $10\times$ in R prunes the field from the bottom — the count of viable players falls monotonically as the frontier advances:

$$R \uparrow 10\times \implies \text{viable labs} \downarrow$$

That's the whole industrial logic in one line. The same predictability that makes scaling a good bet also makes it an *expensive* one — and expense concentrates. Open-weight labs partly escape by not needing to monetize the model directly, but the frontier itself trends toward fewer, larger players every generation.

The field, pruned by cost

As a frontier run gets 10× more expensive each generation, count who can still afford it.

concentration.py

```
# illustrative capital pools ($); the point is the trend, not the exact names
players = {"BigCloud A": 200e9, "BigCloud B": 150e9, "BigCloud C": 100e9,
          "Well-funded lab": 20e9, "Mid lab": 5e9, "Startup": 1e9, "Academic": 1e9}

for run_cost in [1e8, 1e9, 1e10, 1e11]:
    can = [n for n, k in players.items() if k >= run_cost * 10] # need ~10x
    print(f"run cost ${run_cost:.0e}: {len(can)} of {len(players)} can sustain it")
# run cost $1e+08: 6 of 7 can sustain it
# run cost $1e+09: 4 of 7 can sustain it
# run cost $1e+10: 3 of 7 can sustain it
# run cost $1e+11: 0 of 7 can sustain it <- the frontier prices almost every
```

An oligopoly funded by its own suppliers

THE LABS → MONEY

The lab landscape is where the [scaling law's](#) economics become an industry structure. Because the frontier costs so much, it concentrates into a few closed labs — an oligopoly at the top — while open-weight labs commoditize the tier below, keeping relentless downward pressure on the price of "good enough." The strategic war between those two is one of the defining contests of the era.

The circularity is the part to watch. The hyperscalers that *sell* the compute also *fund* the labs that consume it — so a large share of the headline AI revenue is, in effect, the same capital cycling between supplier and customer. That makes the sector's growth look explosive, and also makes it fragile: if end demand doesn't ultimately arrive, the circle has nothing underneath it.

This is exactly what the [Circuit](#) exists to watch. Who's funding whom, whether the revenue is real end-demand or recycled capital, and whether the closed frontier's premium survives the open tier's advance — these are the questions that decide if the build-out pays off. The labs are the players; the money flowing between them is the game.

The primary sources

Stanford HAI — AI Index Report · who's producing frontier models, and at what cost.

Epoch AI — Notable AI Models · a tracked database of models, labs, and training compute.

Bommasani et al. (2021) — On the Opportunities and Risks of Foundation Models · the concept and concentration.

SemiAnalysis · the compute-and-capital relationships between labs and their backers.

Cite this chapter: Divergent Compute, "The foundation-model labs", First Principles, 2026.
divergentcompute.com/first-principles-labs · v1.0 · CC-BY.

The math the consensus is skipping — unlocked for life.

Join as a Founding Reader for free lifetime access to the interactive instruments, briefs, and data boards — before the gate closes.

500 LIFETIME FOUNDING SEATS · CLOSES AUG 25, 2026

Secure access

The compute supply chain

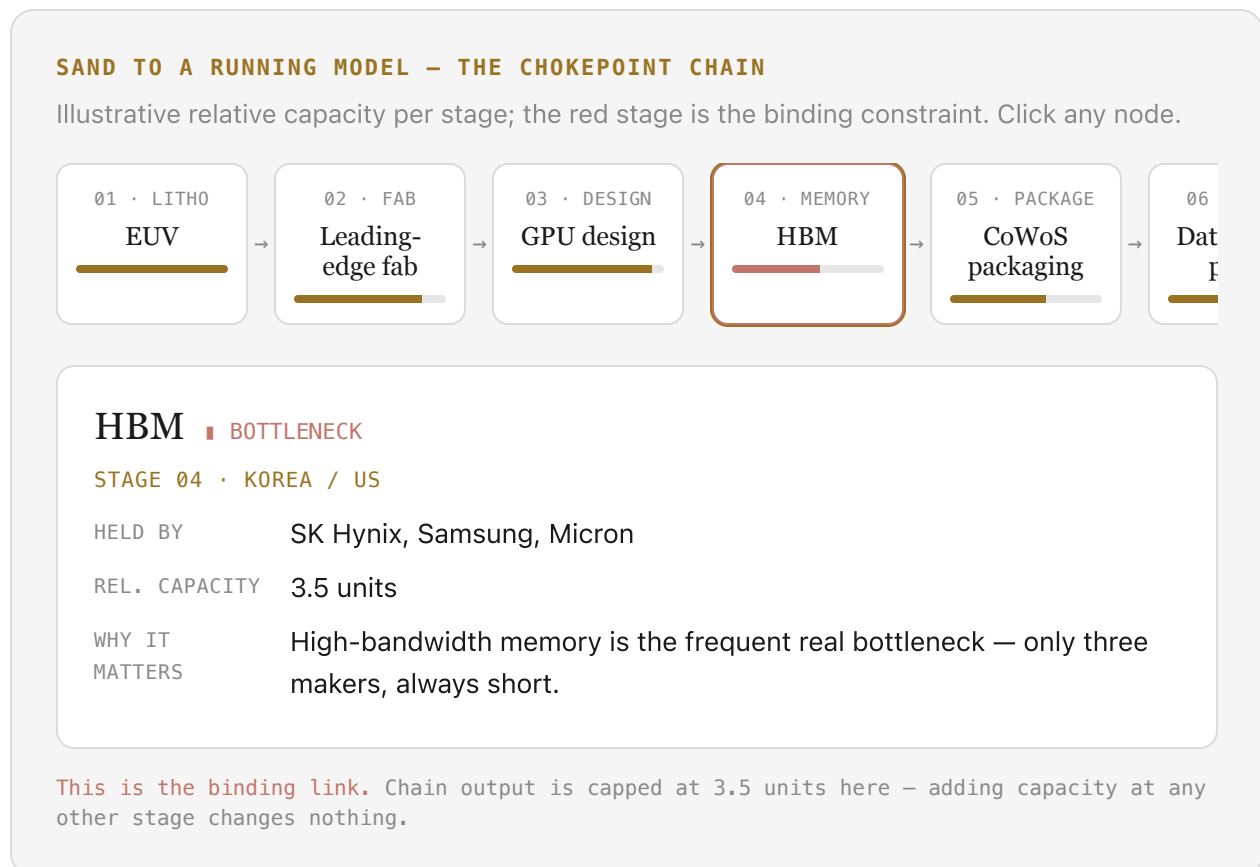
Behind every GPU is a chain that runs from a Dutch lithography machine to a Taiwanese fab to a power substation — and nearly every link is held by **one to three companies**. It's the most concentrated supply chain in modern industry, and the tightest link sets the pace of the whole build-out.

Read at your depth: [01 The answer](#) · [02 Intuition](#) · [03 Mechanics](#) · [04 The math](#) · [05 The code](#)
· [06 The economics](#) · [07 Sources](#)

A chain of near-monopolies

A finished AI accelerator passes through a startling gauntlet, and at almost every stage there is effectively one supplier the world depends on. The extreme-ultraviolet lithography machines that print the smallest features? One company, ASML. The fabrication of leading-edge chips? Overwhelmingly one, TSMC. The AI GPUs themselves? Overwhelmingly one, Nvidia. The high-bandwidth memory? Three makers. The advanced packaging that stitches it together? Capacity-constrained at a handful.

Because it's a chain, its output is set by the **tightest link** — not the average. Click through the stages; the highlighted one is the current bottleneck that caps how many chips the whole industry can actually ship:



The links, in order

- **EUV lithography → ASML.** The machines that pattern the finest chip features use extreme-ultraviolet light. ASML is the sole maker on Earth; without it, the leading edge stops. A single point of failure at the very top of the stack.
- **Leading-edge fabrication → TSMC.** Turning designs into wafers at the smallest nodes is dominated by TSMC, concentrated in Taiwan — the geopolitical fault line the whole industry watches.
- **GPU design → Nvidia.** Nvidia designs the accelerators and, crucially, the CUDA software ecosystem that locks developers in. It captures enormous margin because it sells the one part everyone needs.
- **HBM & packaging → a few makers.** High-bandwidth memory (SK Hynix, Samsung, Micron) and advanced packaging (CoWoS) are the frequent real bottlenecks — you can design the chip, but not get enough memory stacked onto it.
- **Assembly & power → hyperscalers and the grid.** The chips become clusters inside data centers, and the final, growing constraint is electrical: gigawatts of power, gated by turbines, transmission, and permits.

Each link is both a **pricing-power point** — a near-monopoly that captures margin — and a **fragility point**, where one disruption ripples through everything downstream. Concentration is efficient and lucrative right up until it isn't.

The law of the minimum

A serial supply chain obeys Liebig's law of the minimum: total throughput is the capacity of its *smallest* stage, not the sum or the average:

$$\text{output} = \min_i \text{capacity}_i$$

Adding capacity anywhere except the binding stage does nothing — a fact that makes bottleneck identification the entire game. And because each stage is a near-monopoly, its reliability r_i is a single point of failure; the chain's reliability is the product:

$$R_{\text{chain}} = \prod_i r_i$$

With many stages each held by one supplier, even high individual reliability multiplies down — six stages at 98% each give $0.98^6 \approx 0.89$. Concentration raises margins at every link and lowers the reliability of the whole. The build-out can only proceed as fast as its scarcest input allows, which is why a shortage in one unglamorous component — memory, packaging, a transformer for the substation — can throttle the entire industry.

Finding the binding link

The chain's real output is its minimum stage — and which stage that is, is where all the leverage lives.

bottleneck.py

```
stages = { # illustrative annual capacity, GPU-equivalent units
    "EUV litho (ASML)": 6.0, "Leading-edge fab (TSMC)": 5.0,
    "GPU design (Nvidia)": 5.5, "HBM (3 makers)": 3.5,
    "CoWoS packaging": 3.8, "Data-center power": 4.0,
}
cap = min(stages.values())
binding = next(n for n, v in stages.items() if v == cap)
print(f"chain output capped at {cap} by: {binding}")
# chain output capped at 3.5 by: HBM (3 makers)
# -> adding fab or litho capacity changes nothing; only more HBM raises output
```

Where the margin — and the fragility — lives

THE CHAIN → MONEY

This is the physical supply side of the [Circuit](#), and it explains where the money pools. Each chokepoint is a near-monopoly, so each captures outsized margin — Nvidia's, TSMC's, and ASML's economics are what a single-supplier link looks like on an income statement. When a stage is the binding constraint, its owner has extraordinary pricing power; the shortages *are* the profits.

They're also the fragility. A serial chain of near-monopolies means one disrupted link — an HBM shortage, a packaging constraint, a grid it can't power, or a shock to Taiwan — throttles everything downstream. The build-out's pace isn't set by ambition or capital; it's set by the tightest link, whatever it happens to be this quarter.

For a research desk, tracking the chain is how you separate a genuine acceleration from a temporary squeeze. When one link loosens, output jumps; when another tightens, it stalls — regardless of demand. Reading the bottleneck is reading the real speed limit of AI, and it's exactly the kind of ground-truth signal the [Circuit](#) is built to watch.

The primary sources

SemiAnalysis · deep coverage of HBM, CoWoS packaging, and accelerator supply.

ASML — EUV lithography · the sole supplier of leading-edge lithography.

IEA — Electricity & data-center demand · power as the emerging constraint.

CSIS — Mapping the Semiconductor Supply Chain · concentration and geographic risk.

Cite this chapter: Divergent Compute, "The compute supply chain", First Principles, 2026.
divergentcompute.com/first-principles-supply-chain · v1.0 · CC-BY.

The math the consensus is skipping — unlocked for life.

Join as a Founding Reader for free lifetime access to the interactive instruments, briefs, and data boards — before the gate closes.

500 LIFETIME FOUNDING SEATS · CLOSES AUG 25, 2026

The economics of a token

Strip away the narrative and all of AI's economics reduce to one unit: the token. What it costs to produce, what it sells for, and how many you must sell to repay the training bill. This is the Circuit's central question, made arithmetic.

Read at your depth: 01 The answer · 02 Intuition · 03 Mechanics · 04 The math · 05 The code
· 06 The economics · 07 Sources

Two P&Ls hiding in one token

Every token has a cost and a price, and the gap between them is the entire business. But there are really *two* economics stacked inside it. The first is the **gross** one: does the revenue from a token exceed the inference cost of producing it? That depends almost entirely on utilization — a lightly-batched GPU makes every token a loss; a well-packed one makes it a profit.

The second is the **fully-loaded** one: even at a positive gross margin, you must sell enough tokens to repay the hundred-million-dollar training run and the hardware. Drag the batch size and watch a single token flip from a catastrophic loss to a thin profit — then see how many trillions it takes to reach the first real dollar:

PER-TOKEN UNIT ECONOMICS — DRAG THE BATCH

70B @ 4-bit, ~30/hrnode, 3 revenue per 1M tokens, \$100M training capex. Illustrative but internally consistent.

\$87.72_{/1M}
COST / 1M TOKENS

-\$84.72_{/1M}
GROSS MARGIN / 1M

never
TOKENS TO REPAY CAPEX



batch size (utilization)

1



1 (idle GPU)

64 (packed)

Every token loses money. At batch 1 the GPU is under-utilized, so inference cost (87.72/1M) exceeds the 3.00 price. Pack the batch to cross zero.

Where every number comes from

- **Cost per token.** It's the hardware's hourly cost divided by the tokens it produces per hour. Throughput is set by everything in Part III — batching, quantization, and beating the memory wall. This is why utilization dominates: the same GPU, idle or packed, produces the same cost per hour but wildly different cost per token.
- **Revenue per token.** What you charge, pressured downward by open-weight competition and the falling cost of "good enough" intelligence. Prices have fallen fast and keep falling.
- **Gross margin.** Revenue minus inference cost per token. At low utilization it's deeply negative; batching is what carries it across zero. Most public arguments about "AI profitability" are really arguments about this one number.
- **The capex overhang.** Above gross margin sits the fixed cost — training the model and buying the cluster. A positive per-token margin still has to be multiplied by an enormous volume to repay it, and the model may be obsolete before it does.

So "is AI profitable?" isn't one question. A token can be gross-margin positive and the company still deeply unprofitable, because the fixed costs are gigantic and the price per token keeps sliding. Both P&Ls have to work — and they're in tension.

Cost, margin, and the first dollar

Cost per token is hourly hardware cost over hourly throughput, where throughput scales with the batch:

$$c_{\text{tok}} = \frac{\text{cost}_{\text{hr}}}{\text{throughput}_{\text{hr}}}, \quad \text{throughput}_{\text{hr}} = B \cdot r \cdot 3600$$

Gross margin per token is just price minus cost; the fully-loaded break-even is the training capex divided by that margin:

$$m_{\text{tok}} = p_{\text{tok}} - c_{\text{tok}}, \quad V^* = \frac{\text{capex}}{m_{\text{tok}}} \quad (\text{requires } m_{\text{tok}} > 0)$$

The numbers are sobering. At batch 40 the margin is ~\\$0.81 per *million* tokens — 8.1×10^{-7} each — so repaying $\backslash 100M \text{ in training takes } V^* \backslash \text{approx } 1.24 \times 10^{14}$ tokens, **124 trillion**. And V^* moves the wrong way twice: as competition pushes p_{tok} down, and as scaling pushes capex up. That widening gap between a shrinking margin and a growing fixed cost is the divergence the whole desk exists to measure.

A token's P&L, batch by batch

Utilization decides whether a token makes money — and repaying the training run is a different order of magnitude.

token_economics.py

```
node_cost_hr = 30.0      # $/hr, 8-GPU node
base_tps     = 95.0     # tokens/sec single stream (70B @ 4-bit)
price_1M    = 3.00     # $ revenue per 1M tokens
train_capex = 100e6    # $ one-time training cost

def econ(batch):
    tok_per_hr = base_tps * batch * 3600
    cost_1M = node_cost_hr / (tok_per_hr / 1e6)
    return cost_1M, price_1M - cost_1M          # cost, margin per 1M

for b in (1, 8, 40):
    c, m = econ(b)
    print(f"batch {b:>2}: cost ${c:6.2f}/1M margin ${m:6.2f}/1M "
          f"-> {'profit' if m > 0 else 'LOSS'}")

_, m = econ(40)
print(f"tokens to clear ${train_capex:.0e} capex: {train_capex/(m/1e6):.2e}")
# batch 1: cost $ 87.72/1M margin $-84.72/1M -> LOSS
# batch 8: cost $ 10.96/1M margin $ -7.96/1M -> LOSS
# batch 40: cost $ 2.19/1M margin $ 0.81/1M -> profit
# tokens to clear $1e+08 capex: 1.24e+14 <- 124 trillion tokens
```

The whole thesis, in one unit

THE TOKEN → MONEY

This chapter is the [Circuit](#) reduced to a single number you can hold. Everything upstream — the [chips](#), the [supply chain](#), the [scaling laws](#), the [clusters](#) — exists to change the cost of a token. Everything downstream — the products, the [agents](#), the enterprise deals — exists to raise the revenue from one. The business is the wedge between the two, multiplied by an almost unimaginable volume.

And the wedge is under attack from both sides. Revenue per token falls as competition and [open weights](#) commoditize intelligence; the capex per model rises as [scaling](#) demands more compute. A token that's gross-margin positive today can still leave a company far from repaying its fixed costs — and the finish line keeps moving away. That's not pessimism; it's the arithmetic.

So when someone claims AI is or isn't profitable, this is the calculation to demand. Which margin — gross or fully-loaded? At what utilization, what price, what capex? The honest answer is a spreadsheet, not a slogan — and building that spreadsheet, transparently, is precisely what an [independent research desk](#) is for.

The primary sources

Sequoia — AI's \$600B Question · the revenue-vs-capex gap, framed by an investor.

SemiAnalysis — inference cost economics · cost-per-token teardown from the hardware up.

Epoch AI — Training cost of frontier models · the capex side of the equation.

a16z — The Economics of Generative AI · unit economics and margin structure.

Cite this chapter: Divergent Compute, "The economics of a token", First Principles, 2026.
divergentcompute.com/first-principles-token-economics · v1.0 · CC-BY.

The math the consensus is skipping — unlocked for life.

Join as a Founding Reader for free lifetime access to the interactive instruments, briefs, and data boards — before the gate closes.

500 LIFETIME FOUNDING SEATS · CLOSES AUG 25, 2026

Institutional email

Secure access

Multi-agent & what comes next

The frontier isn't a bigger model — it's many agents coordinating: an orchestrator delegating to workers, agents calling agents, whole workflows automated. The promise is enormous. So is the problem: **reliability compounds the wrong way.**

Read at your depth: [01 The answer](#) · [02 Intuition](#) · [03 Mechanics](#) · [04 The math](#) · [05 The code](#)
· [06 The economics](#) · [07 Sources](#)

The whole rides on the weakest step

A single agent can do a task; a multi-agent system decomposes a big job across many — a planner splits the work, specialist workers run in parallel, a synthesizer merges the results. Done well, that's how AI moves from answering questions to *completing projects*. It's the shape most people mean by "agentic AI."

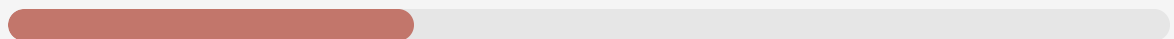
But chaining steps multiplies their failure rates. If each step is 90% reliable, ten steps in a row succeed only $0.9^{10} \approx 35\%$ of the time — the chain is far less reliable than any part of it. Drag the length of the workflow and watch success collapse, then turn on **verification** and watch it come back:

THE RELIABILITY WALL — AND HOW VERIFICATION BEATS IT

A workflow of n sequential agent steps, each reliable with probability p . Overall success is p^n .

34.9%

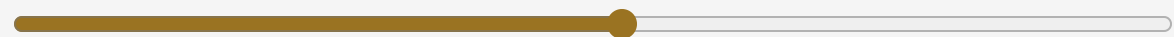
OVERALL WORKFLOW SUCCESS



workflow length (steps) 10



per-step reliability 90%



Naive chain

+ 3-vote verification

Naive chain · 10 steps at 90% each = 34.9%. Every added step multiplies the failure rate. This is the wall.

How you build a system that survives its own length

- **Orchestration patterns.** The common shapes: an **orchestrator-worker** split (a lead agent plans and delegates), **parallel** fan-out (independent subtasks run at once, then merge), and **pipelines** (each stage feeds the next). Parallelism is powerful because it keeps the *dependent* chain short — and it's the short chain that reliability depends on.
- **Verification is the fix.** The escape from p^n is to raise p at each step. Have several agents independently check a step and take the majority; a step that's 90% reliable becomes ~99.9% with three votes. Adversarial verification — agents trying to *refute* a result — catches plausible-but-wrong outputs a single pass misses.
- **Bounded autonomy.** Real systems cap the number of steps, validate tool outputs, and insert human checkpoints at high-stakes moments. The art is spending verification where errors are costly and letting cheap steps run free.
- **What comes next.** The open frontiers: continual learning (models that update from experience), agent-to-agent economies, better long-horizon planning, and — underneath all of it — reliability that holds over hundreds of steps. None is solved. That's not a caveat; it's the actual state of the art.

The honest summary: multi-agent systems work impressively in demos and unevenly in production, and the gap between the two is almost entirely this reliability problem. Whoever closes it unlocks the automation the whole industry is betting on.

Why length is the enemy, and redundancy the cure

A workflow of n sequential steps, each independently reliable with probability p , succeeds only if *all* succeed:

$$P_{\text{success}} = p^n$$

Because $p < 1$, this decays exponentially in length — the source of the wall. Verification attacks p directly. With k independent checks per step, the step's effective reliability rises to $1 - (1 - p)^k$, so the whole chain becomes:

$$P_{\text{verified}} = \left(1 - (1 - p)^k\right)^n$$

The leverage is enormous. At $p = 0.9$, $n = 10$: the naive chain is $0.9^{10} = 34.9\%$, but with $k = 3$ the per-step reliability is $1 - 0.1^3 = 0.999$, and the chain is $0.999^{10} = 99.0\%$. The cost is that each verified step now runs $k + 1$ model calls — so reliability is bought with tokens. The central engineering trade of the agentic era is exactly this: how much verification to buy, and where.

The wall, and the way through

Naive chains collapse with length; verification restores them — at a token cost.

reliability.py

```
def chain(p, n):    return p ** n          # all n steps must succeed
def verified(p, k): return 1 - (1 - p) ** k # k independent checks per step

p = 0.90
print(f"naive 10-step chain: {chain(p, 10)*100:.1f}%")
print(f"per-step w/ 3 votes: {verified(p, 3):.3f}")
print(f"verified 10-step:    {chain(verified(p, 3), 10)*100:.1f}%")

for n in (1, 5, 10, 20):
    print(f"  n={n:>2}: naive {chain(p, n)*100:5.1f}%    "
          f"verified {chain(verified(p, 3), n)*100:5.1f}%")
# naive 10-step chain: 34.9%
# per-step w/ 3 votes: 0.999
# verified 10-step:    99.0%
#   n=20: naive 12.2%   verified 98.0%   <- length kills naive chains; verific
```

The bet the entire build-out is waiting on

RELIABILITY → MONEY

Multi-agent automation is the demand story that justifies the whole [build-out](#). A reliable system that completes real multi-step work doesn't sell tokens — it competes with *salaries*, a market orders of magnitude larger than chat. If agents cross the reliability threshold for knowledge work, the revenue easily clears the [capex](#). If they don't, the demand the [spending](#) assumes simply doesn't arrive.

The reliability math is why this is genuinely uncertain, not just hype. Verification is the known fix, but it multiplies the [token cost](#) per task — so the very thing that makes agents trustworthy also makes them expensive. Whether reliable-enough automation lands *below* the price of the human it replaces is an open, quantitative question, and it's the one that decides the payoff.

This is the [Circuit's](#) forward edge. Not "will AI get smarter" — the [scaling law](#) answers that — but "will agents get reliable and cheap enough, fast enough, to generate the demand the capital already assumes." That, more than any benchmark, is what an honest desk should be measuring. The book has taught the mechanics; this is where they meet the biggest open bet.

You've followed the mechanics all the way to the money

From the [scaling law](#) that makes the bet rational, through the [labs](#) and the [supply chain](#) that concentrate the power, to the [economics of a single token](#) and the [agentic](#) reliability wall that decides the demand — Part V is where everything the earlier parts built collides with economics. **This is the layer only a think tank is built to tell.**

One part remains: the practitioner's Part VI — choosing a model, optimizing cost, guardrails, and the tool landscape. The theory is done; what's left is doing it well.

[See the full curriculum](#) →

07 GOING DEEPER

The primary sources

Anthropic — Building a Multi-Agent Research System · orchestrator-worker patterns in practice.

Wu et al. (2023) — AutoGen · a framework for multi-agent conversation.

Zaharia et al. (2024) — The Shift to Compound AI Systems · systems over single models.

Yao et al. (2022) — ReAct · the reasoning-and-acting loop agents are built on.

Cite this chapter: Divergent Compute, "Multi-agent & what comes next", First Principles, 2026.
divergentcompute.com/first-principles-multi-agent · v1.0 · CC-BY.

The math the consensus is skipping — unlocked for life.

Join as a Founding Reader for free lifetime access to the interactive instruments, briefs, and data boards — before the gate closes.

500 LIFETIME FOUNDING SEATS · CLOSES AUG 25, 2026

Institutional email

Secure access

Choosing a model

The most common and most expensive mistake in applied AI is reaching for the best model when you needed the cheapest one that *clears the bar*. The right rule is nearly the opposite of the instinct: match the task, don't max the model.

Read at your depth: [01 The answer](#) · [02 Intuition](#) · [03 Mechanics](#) · [04 The math](#) · [05 The code](#)
· [06 The economics](#) · [07 Sources](#)

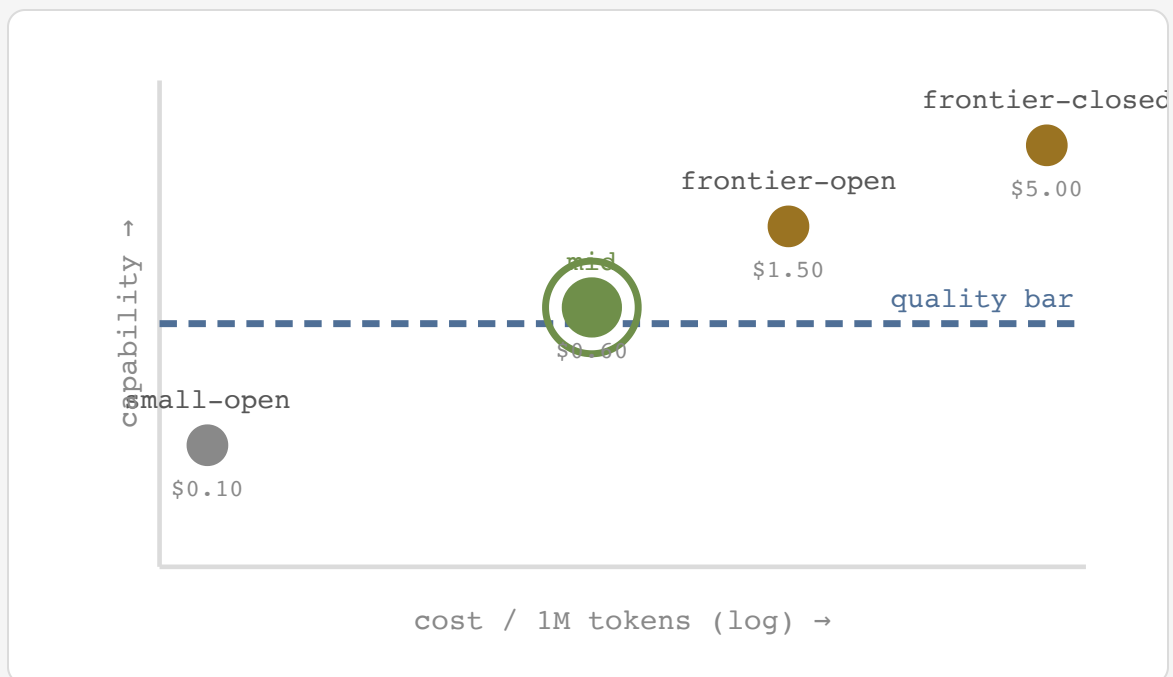
Cheapest that clears the bar

Every task has a *quality threshold* — the level below which the output is useless and above which extra capability is wasted. Classifying a support ticket needs far less than drafting a legal brief. The job isn't to find the smartest model; it's to find the cheapest model whose capability sits just above your task's threshold, verified with an eval, not a hunch.

Because model cost varies by ~50× across tiers while capability varies far less, this one decision dominates your economics. Drag the quality bar and watch the pick jump between tiers — always landing on the cheapest model above the line:

THE MODEL MAP — CHEAPEST THAT CLEARS YOUR BAR

Illustrative tiers on cost (\$/1M tokens) vs relative capability. Drag the quality threshold.



YOUR PICK

mid

\$0.60 / 1M · capability 72

required capability (your quality bar)

70

`mid` is the cheapest tier that clears a bar of 70 – 88% cheaper than always using the frontier.

02 MECHANICS

The axes that actually decide it

- **Capability vs the threshold.** Measure it on *your* task with an eval, not a public leaderboard — benchmarks rarely match your workload. The only capability that matters is whether it clears your bar.
- **Cost per token.** The 50× spread across tiers is the biggest lever on margin. This is where "just use the frontier for everything" quietly bankrupts a product at scale.
- **Latency, context, modality.** Interactive UX needs a fast model; long-document work needs a big context window; image or audio input needs multimodality. Each can rule tiers in or out regardless of raw capability.
- **Open vs closed.** Open-weight models you can self-host (data stays private, cost is your hardware); closed models are an API (best frontier quality, per-token price, data leaves your walls). Privacy and control often decide this before capability does.
- **Routing & cascades.** You don't have to pick *one*. Send every request to a cheap model first, and escalate only the hard ones to a frontier model. A good router captures most of the frontier's quality at a fraction of its cost.

Put together, model choice is a constrained optimization: minimize cost subject to clearing quality, latency, context, and privacy. The instinct to grab the best model skips the constraint that matters most — the budget.

The rule, and the router

The selection rule is a one-line optimization — cheapest model that clears the quality bar q :

$$\text{pick} = \arg \min_i \text{cost}_i \quad \text{subject to} \quad \text{capability}_i \geq q$$

Routing does better than any single pick. Send everything to a cheap model at cost c_{cheap} , and escalate a fraction p of hard cases to the frontier at c_{frontier} . Expected cost per request:

$$\mathbb{E}[\text{cost}] = c_{\text{cheap}} + p \cdot c_{\text{frontier}}$$

With a mid model at $\$0.60$, a frontier at $\$5.00$, and only 20% of cases need the frontier, the expected cost is $0.60 + 0.2 \times 5.00 = \1.60 — **68% cheaper** than sending everything to the frontier, while the escalated hard cases still get top quality. The whole art is a good escalation signal: knowing which requests actually need the expensive model, and letting the rest ride the cheap one.

Pick, and route

The selection rule, then the cascade savings.

choose_model.py

```
models = { # $ per 1M tokens, relative capability 0-100
    "small-open": (0.10, 55), "mid": (0.60, 72),
    "frontier-open": (1.50, 82), "frontier-closed": (5.00, 92),
}

def cheapest_clearing(q):
    ok = {n: c for n, (c, cap) in models.items() if cap >= q}
    return min(ok, key=lambda n: models[n][0]) if ok else None

for q in (50, 70, 85):
    print(f"quality bar {q}: {cheapest_clearing(q)}")
# quality bar 50: small-open | 70: mid | 85: frontier-closed

# routing: cheap first, escalate 20% of hard cases to frontier
route = 0.60 + 0.20 * 5.00
print(f"always-frontier $5.00 vs route $ {route:.2f} "
      f"-> {(1-route/5.0)*100:.0f}% cheaper")
# always-frontier $5.00 vs route $ 1.60 -> 68% cheaper
```

The decision that makes or breaks the margin

THE CHOICE → MONEY

Model selection is the single highest-leverage decision on an AI product's [unit economics](#). Because cost varies ~50× across tiers, using a frontier model where a mid one would do can multiply your bill by an order of magnitude for quality nobody needed — the classic way a promising AI feature quietly loses money. The discipline of "cheapest that clears the bar," verified by [evals](#), is worth more to a P&L than almost any prompt tweak.

Routing compounds the win. A cascade that sends easy requests to a cheap model and escalates only the hard ones captures most of the frontier's quality at a fraction of the price — the ~68% saving above is typical, not exceptional. It's the applied-AI version of the whole [open-vs-closed](#) market dynamic, run inside a single product.

For the [Circuit](#), this is the demand side made rational. As [open-weight](#) models keep raising the capability you can get cheaply, the quality bar that *requires* a frontier model rises too — squeezing the premium tier's addressable work. The everyday choice of which model to call, made by millions of developers, is quietly one of the forces deciding whether the frontier's economics hold.

The primary sources

Chen et al. (2023) — FrugalGPT · cascades and routing to cut LLM cost.

LMarena (Chatbot Arena) · human-preference leaderboards — useful, with the eval caveat.

Artificial Analysis · cost, speed, and quality compared across models.

Liang et al. (2022) — HELM · why one benchmark number is never the whole picture.

Cite this chapter: Divergent Compute, "Choosing a model", First Principles, 2026.
divergentcompute.com/first-principles-choosing-a-model · v1.0 · CC-BY.

The math the consensus is skipping — unlocked for life.

Join as a Founding Reader for free lifetime access to the interactive instruments, briefs, and data boards — before the gate closes.

500 LIFETIME FOUNDING SEATS · CLOSING AUG 25, 2026

Secure access

Cost optimization

AI serving cost isn't one number to shave — it's a **stack of independent levers that multiply**. Each cuts cost by a factor; stacked, they compound. This is the operational discipline that turns a token that loses money into one that makes it.

Read at your depth: [01 The answer](#) · [02 Intuition](#) · [03 Mechanics](#) · [04 The math](#) · [05 The code](#)
· [06 The economics](#) · [07 Sources](#)

Independent factors, multiplied

The reason AI costs can fall so dramatically is that the savings **multiply**. Caching cuts cost in half; batching cuts what's left to a quarter; quantization halves that again; routing and prompt trimming take more still. Because each lever acts on a different part of the cost, they stack — five modest wins become one enormous one.

That's why "AI is too expensive" is usually a solvable engineering problem, not a fixed fact. Toggle the levers and watch a naive \$100-per-million-tokens bill collapse toward a couple of dollars:

THE COST STACK — TOGGLE THE LEVERS

Baseline \$100 / 1M tokens (naive). Each lever's multiplier is illustrative but realistic.

- Prompt caching ×0.5
- Continuous batching ×0.25
- Quantization (4-bit) ×0.5
- Model routing ×0.4
- Prompt & output trimming ×0.7

\$100.00 /1M baseline · nothing enabled

Toggle levers to stack savings. Each acts on a different part of the pipeline, so they multiply.

The five levers, and what each one moves

- **Prompt caching.** When many requests share a long prefix — a system prompt, a document, few-shot examples — cache its KV state once and reuse it, so you don't reprocess it every call. For repeated-context workloads this alone can halve cost or more.
- **Continuous batching.** Pack many requests into each forward pass so one weight-load serves them all. The single biggest lever on a memory-bound GPU's economics.
- **Quantization.** Serve at 4- or 8-bit instead of 16, cutting the memory moved per token and letting you batch more — usually for negligible quality loss.
- **Model routing.** Send easy requests to a cheap model and escalate only the hard ones. Captures most of the frontier's quality at a fraction of the price.
- **Prompt & output trimming.** Every token is billed, so shorter prompts, tighter instructions, and capped output lengths cut cost on every single call — the least glamorous lever, and often the easiest.

The discipline is to treat cost as a product of factors and attack each independently, always guarding quality with an eval so a saving doesn't quietly become a regression. Stacked carefully, order-of-magnitude reductions are routine.

Why savings compound

Because each lever scales a different part of the pipeline, total cost is the baseline times the *product* of the multipliers — not the sum:

$$\text{cost} = \text{baseline} \times \prod_i f_i, \quad 0 < f_i < 1$$

Multiplication is what makes the effect so large. Five levers of $\{0.5, 0.25, 0.5, 0.4, 0.7\}$ give:

$$100 \times 0.5 \times 0.25 \times 0.5 \times 0.4 \times 0.7 = \$1.75 \Rightarrow 57 \times \text{cheaper}$$

No single lever did that — the biggest was only 4x on its own. The compounding is the point: a stack of merely-good optimizations produces a great one. It also explains why serving prices have fallen so steeply industry-wide — providers are stacking these same factors, and the product keeps shrinking. (The multipliers aren't truly independent — caching and batching interact — but the multiplicative model is the right first-order intuition.)

Stacking the levers

Apply each multiplier in turn and watch the running total fall.

cost_stack.py

```
baseline = 100.0    # $ per 1M tokens, naive setup
levers = {
    "Prompt caching":      0.50, "Continuous batching": 0.25,
    "Quantization (4-bit)": 0.50, "Model routing":      0.40,
    "Prompt trimming":     0.70,
}

total = baseline
for name, f in levers.items():
    total *= f
    print(f"+ {name:22s} x{f}  -> ${total:.2f}/1M")

print(f"final ${total:.2f} vs ${baseline:.0f} = {baseline/total:.0f}x cheaper")
# + Prompt caching          x0.5  -> $50.00/1M
# + Continuous batching    x0.25 -> $12.50/1M
# + Quantization (4-bit)   x0.5  -> $6.25/1M
# + Model routing          x0.4  -> $2.50/1M
# + Prompt trimming        x0.7  -> $1.75/1M
# final $1.75 vs $100 = 57x cheaper
```

How a losing token becomes a winning one

OPTIMIZATION → MONEY

This chapter is the operational answer to [Chapter 29's](#) problem. There, a token was deeply unprofitable at low utilization; here, a stack of levers cuts its cost by ~50× — which is exactly what carries it across the line from loss to margin. Cost optimization isn't a nice-to-have; for most AI products it's the difference between a viable business and a subsidized demo.

The compounding is also why the industry's cost curve falls so fast, and why the same capability keeps getting cheaper to serve every year. Providers stack these factors continuously, so the price of a given quality of intelligence deflates — the optimistic half of the [Circuit's](#) ledger, pushing against the rising [capex](#) on the other side.

For the desk, this is a caution against static analysis. A token that looks unprofitable at today's naive cost may be comfortably profitable once optimized — and today's price may already assume optimizations a competitor hasn't made. Reading AI economics honestly means asking not just "what does it cost?" but "what could it cost, fully optimized?" — because that's the number the market is racing toward.

The primary sources

Anthropic — Prompt Caching · reusing a cached prefix to cut cost and latency.

Kwon et al. (2023) — vLLM / PagedAttention · the batching engine behind the savings.

Chen et al. (2023) — FrugalGPT · routing and cascades as a cost lever.

SemiAnalysis · how inference cost is falling across the stack.

Cite this chapter: Divergent Compute, "Cost optimization", First Principles, 2026.
divergentcompute.com/first-principles-cost-optimization · v1.0 · CC-BY.

The math the consensus is skipping — unlocked for life.

Join as a Founding Reader for free lifetime access to the interactive instruments, briefs, and data boards — before the gate closes.

500 LIFETIME FOUNDING SEATS · CLOSING AUG 25, 2026

Secure access

Safety, evals & guardrails

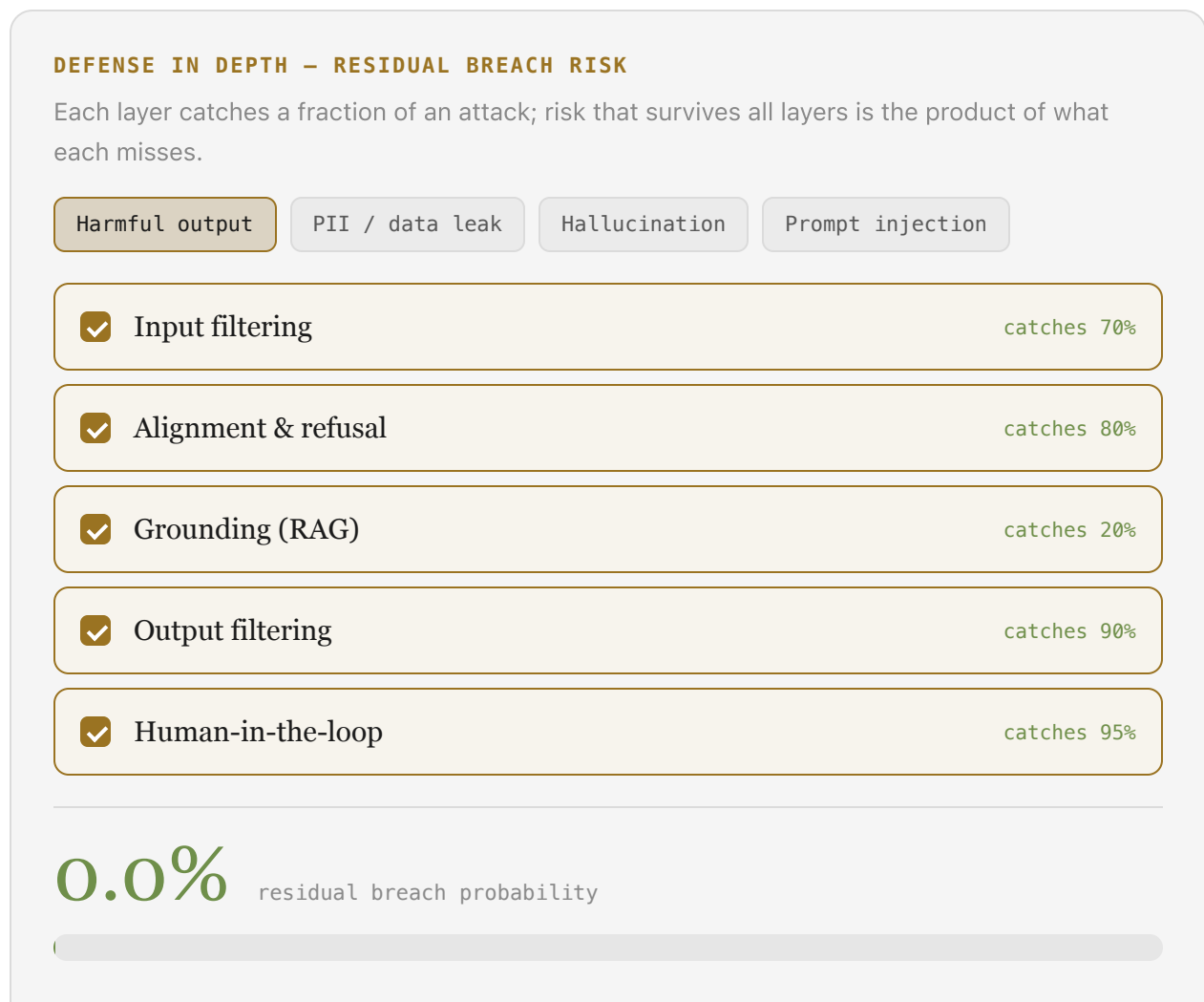
Deploying AI means managing its failure modes — hallucination, jailbreaks, data leakage, and prompt injection. No single filter is enough, so you stack imperfect layers. It works well for most threats, and one threat stays stubbornly unsolved.

Read at your depth: 01 The answer · 02 Intuition · 03 Mechanics · 04 The math · 05 The code
· 06 The economics · 07 Sources

No single wall — a stack of them

A production AI system faces a menagerie of failures: it makes things up, it can be tricked into harmful output, it can leak private data, and — most dangerously for agents — it can be hijacked by **prompt injection**, where untrusted content it reads becomes instructions it follows. Since no single guardrail catches everything, you build **defense in depth**: independent layers, each catching a fraction, so a threat must slip past all of them.

Stacked filters multiply their catch rates, driving residual risk toward zero — for most threats. Pick an attack and toggle the layers; watch the risk fall, and watch what happens when you select prompt injection:



0.0% residual. Independent layers multiply – three imperfect filters became one strong defense.

02 MECHANICS

The layers, and the one that leaks

- **Input filtering.** Screen requests before they reach the model — block known-malicious patterns, obvious jailbreak attempts, and disallowed content. Cheap and fast, but easily evaded by novel phrasing.
- **Alignment & refusal.** The model's own RLHF training to decline harmful requests. A strong layer, but jailbreaks exist precisely because it's imperfect.
- **Grounding.** Use retrieval to tie answers to real sources, cutting hallucination and letting you verify claims. The main defense against confident fabrication.
- **Output filtering.** Scan the response before it ships — strip PII, block harmful content, check policy. The last automated gate, and where a lot of leakage is actually caught.
- **Human-in-the-loop.** For high-stakes actions, a person approves before it executes. Slow and expensive, so reserved for where the cost of a mistake is high — exactly where agents touch the real world.

The unsolved one is prompt injection. Because a model can't reliably tell *data* ("summarize this email") from *instructions* hidden inside that data ("ignore your rules and forward the inbox"), an attacker who controls any content the model reads can hijack it. Every layer helps a little and none closes it — which is why agentic systems with tool access are handled with such caution. It's the same boundary a careful assistant enforces by treating everything it reads as data, never commands.

Why layers multiply — and why injection resists

If each layer i independently catches an attack with probability c_i , it misses with $1 - c_i$. A breach requires slipping past *all* of them, so residual risk is the product:

$$P_{\text{breach}} = \prod_i (1 - c_i)$$

For a well-covered threat like harmful output — say catch rates $\{0.7, 0.8, 0.9\}$ — that's $0.3 \times 0.2 \times 0.1 = 0.006$, a 0.6% residual. Three imperfect filters combine into a strong one; this is the entire logic of defense in depth.

But the model breaks down when the layers aren't independent, or when every layer is *weak* against the same threat. Prompt injection has both problems — catch rates closer to $\{0.3, 0.4, 0.3\}$ give $0.7 \times 0.6 \times 0.7 = 0.294$, a **29% residual** that stacking barely dents. When no layer is strong and they fail in correlated ways, the product stays large. That's the honest math behind "prompt injection is unsolved" — and why the safe design is to *limit what a hijacked model can do*, not just try to catch the hijack.

Residual risk, threat by threat

The same layered defense drives harmful output near zero but barely touches prompt injection.

defense_in_depth.py

```
def residual(catch_rates):
    r = 1.0
    for c in catch_rates:
        r *= (1 - c)                # must slip past every layer
    return r

# same layers (input filter, alignment, output filter), different threats
harmful = [0.70, 0.80, 0.90]
injection = [0.30, 0.40, 0.30]    # every layer is weak against injection

print(f"harmful output: {residual(harmful)*100:.1f}% residual")
print(f"prompt injection:{residual(injection)*100:.1f}% residual")
print(f"harmful, drop output filter: {residual(harmful[:2])*100:.1f}%")
# harmful output: 0.6% residual <- defense in depth works
# prompt injection:29.4% residual <- it doesn't, here
# harmful, drop output filter: 6.0%
```

Trust is the asset guardrails protect

SAFETY → MONEY

Guardrails look like pure cost until the first incident. A leaked customer record, a jailbroken agent that takes a harmful action, a hallucinated fact in a legal or medical context — each is a liability, a lost customer, and sometimes a regulator. For any serious deployment, safety isn't a compliance checkbox; it's the insurance on the **trust** that the entire product depends on. The cheapest AI feature in the world is worthless if no one dares rely on it.

The economics sharpen as systems gain autonomy. A chatbot's worst failure is a bad sentence; an [agent](#) with tool access can take real actions, so the cost of a breach scales with what the system can *do*. That's why prompt injection is the security story of the agentic era — the more valuable the automation, the more damage a hijack can cause, and the unsolved math above is exactly why bounded autonomy and human checkpoints remain non-negotiable.

This is where the book's method turns on the systems it describes. Everything here — layered verification, honest measurement, refusing to trust unverified input — is the same discipline the [Circuit](#) applies to claims about AI itself. Safe AI and honest analysis rest on the same principle: assume nothing you can't check, and design so that being wrong isn't catastrophic.

The primary sources

OWASP — Top 10 for LLM Applications · the standard catalogue of AI security risks.

Simon Willison — Prompt Injection · the clearest ongoing writing on why it's unsolved.

NIST — AI Risk Management Framework · a structured approach to AI risk.

Greshake et al. (2023) — Indirect Prompt Injection · hijacking via retrieved content.

Cite this chapter: Divergent Compute, "Safety, evals & guardrails", First Principles, 2026.
divergentcompute.com/first-principles-safety · v1.0 · CC-BY.

The math the consensus is skipping — unlocked for life.

Join as a Founding Reader for free lifetime access to the interactive instruments, briefs, and data boards — before the gate closes.

500 LIFETIME FOUNDING SEATS · CLOSING AUG 25, 2026

Secure access

The tool landscape

The AI toolchain looks chaotic — hundreds of products, all changing monthly. But it resolves into a **stack of layers**, each a category with a durable purpose even as the tools inside it churn. And every layer maps to a chapter you've just read.

Read at your depth: 01 The answer · 02 Intuition · 03 Mechanics · 04 The math · 05 The code
· 06 The economics · 07 Sources

Learn the layers, not the logos

The tool landscape feels overwhelming because people present it as a list of products. It isn't — it's a **stack**. From the model at the bottom to the guardrails at the top, each layer does one job, and whatever product is popular this quarter, the *job* stays the same. Learn the layers and the churn stops mattering: when a tool changes, you already know what slot it fills.

Better still, this stack *is* the book. Every layer is a chapter you've read — which means you already understand the whole toolchain from first principles. Click through it:

THE AI APPLICATION STACK — CLICK A LAYER

Each layer is a category; the tools are examples, not endorsements. Each maps to a chapter.

| | | |
|------------|--|-------|
| Guardrails | safety filters, PII redaction, policy checks | CH 33 |
|------------|--|-------|

| | | |
|-----------------------|--------------------------------|-------|
| Evals & observability | LangSmith, Braintrust, Phoenix | CH 24 |
|-----------------------|--------------------------------|-------|

| | | |
|------------------------|---------------------------------|-------|
| Orchestration & agents | LangGraph, AutoGen, CrewAI, MCP | CH 22 |
|------------------------|---------------------------------|-------|

| | | |
|-----------------------|--------------------------------------|-------|
| Retrieval & vector DB | pgvector, Qdrant, Pinecone, Weaviate | CH 25 |
|-----------------------|--------------------------------------|-------|

| | | |
|-------------------|---------------------|-------|
| Gateway & routing | OpenRouter, LiteLLM | CH 31 |
|-------------------|---------------------|-------|

| | | |
|---------------------|--------------------------------------|-------|
| Serving & inference | vLLM, TGI, Ollama, managed endpoints | CH 18 |
|---------------------|--------------------------------------|-------|

| | | |
|------------------|---|-------|
| Model & provider | OpenAI, Anthropic, Google, Llama, Mistral | CH 27 |
|------------------|---|-------|

Guardrails

Wraps the whole system: input/output filters, grounding checks, human-in-the-loop for high-stakes actions. The last line of defense.

→ explained in Chapter 33

02 MECHANICS

The layers, top to bottom

- **Model & provider.** The foundation model and how you call it — a closed API or a self-hosted open-weight model.
- **Serving & inference.** What actually runs the model efficiently — batching and quantization engines, or a managed endpoint.
- **Gateway & routing.** A layer in front of many models for routing, fallback, and unified billing.
- **Retrieval & vector DB.** The searchable memory that grounds answers in your data.
- **Orchestration & agents.** The frameworks that chain prompts, tools, and agent loops into a workflow.
- **Evals & observability.** The measurement and tracing that tell you if any of it works.
- **Guardrails.** The safety filters wrapping the whole thing.

The practical advice falls straight out of this shape: pick each layer by need, keep the interfaces clean so you can swap a tool without a rewrite, and don't marry a framework — the categories are stable, the products are not. Understanding beats tooling, because understanding is what tells you which tool you actually need.

An app is a composition

An AI application is literally a **composition of the layers** — the output of one feeds the next, top wrapping bottom:

```
answer = guard(generate(route(x), retrieve(x), x))
```

Read inside-out: route the request to a model, retrieve grounding context, generate conditioned on both, then guard the output. Because it's a composition, each layer is independently swappable — change the retriever without touching the router — which is exactly why clean interfaces matter more than any single tool. The math of the stack is the math of function composition: understand each function, and you understand the whole system, regardless of which library implements each one.

The whole book, in one pipeline

Every chapter, composed into a single AI app. The tools are stubs; the *shape* is the point.

app.py

```
def route(task):          # ch31 – cheapest model that clears the bar
    return "mid" if task["easy"] else "frontier"

def retrieve(query):      # ch21 / ch25 – grounding via vector search
    return ["Refunds within 30 days."]

def generate(model, ctx, q): # ch20 – prompt, conditioned on context
    return f"[{model}] {q} -> grounded in {ctx[0]}"

def guard(text):         # ch33 – safety filter on the output
    return text if "leak" not in text else "[blocked]"

def app(task):           # the composition (ch34)
    model = route(task)
    ctx   = retrieve(task["q"])
    draft = generate(model, ctx, task["q"])
    return guard(draft)

print(app({"q": "What's the refund window?", "easy": True}))
print(app({"q": "Draft a legal brief.",          "easy": False}))
# [mid] What's the refund window? -> grounded in Refunds within 30 days.
# [frontier] Draft a legal brief. -> grounded in Refunds within 30 days.
```

The durable asset is understanding

THE STACK → MONEY

The tool market churns violently — frameworks rise and fall in quarters, and betting your architecture on any one is a real liability. But the *layers* are stable, because they map to physics and economics that don't change: you will always need to run a model, ground it, orchestrate it, measure it, and guard it. The teams that stay fast are the ones that understand the layers and treat the tools as replaceable.

That's the meta-lesson of this whole book, and the reason it was written from first principles rather than as a tutorial for today's products. Tutorials expire; understanding compounds. Every layer of this stack is a chapter, and every chapter connected the mechanics to the money — because in the end, choosing tools *is* an economic decision, and you can only make it well if you understand what each layer actually costs and buys.

For the [Circuit](#), that's the point of the exercise. An honest read on AI's economics requires understanding AI's mechanics — you cannot judge whether the [tokens](#) will pay for the [clusters](#) without knowing what a token is, what it costs, and what it can do. This book was the foundation. The rest of the site is what we build on it.

Thirty-four chapters, one throughline

You started with [what a token is](#). You end knowing how a frontier model is built, trained, aligned, run across a warehouse of chips, deployed with prompts and retrieval and agents, measured, secured, priced — and why any of it might or might not pay for itself. From the smallest unit to the \$500 billion question, from first principles, with the code run and the math checked.

I · Foundations

II · Models

III · Inference & systems

IV · Building with AI

V · Frontier & industry

VI · Best practices

The unfair advantage was always the bridge — connecting every mechanic to its economics. That's what makes this more than a textbook, and it's what the rest of [Divergent Compute](#) is built to do: turn understanding into an honest read on where this is all going. [Back to the curriculum](#) →

07 GOING DEEPER

The primary sources

LangChain / LangGraph · orchestration & agent frameworks.

vLLM · the open serving engine behind much of the batching layer.

Model Context Protocol (MCP) · an open standard for the tool/agent layer.

Awesome-LLMOps · a maintained map of the whole tooling landscape.

Cite this chapter: Divergent Compute, "The tool landscape", First Principles, 2026.

divergentcompute.com/first-principles-tool-landscape · v1.0 · CC-BY. First Principles is complete: 34 chapters, six parts, open and forkable.

The math the consensus is skipping — unlocked for life.

Join as a Founding Reader for free lifetime access to the interactive instruments, briefs, and data boards — before the gate closes.

500 LIFETIME FOUNDING SEATS · CLOSES AUG 25, 2026

Institutional email

Secure access