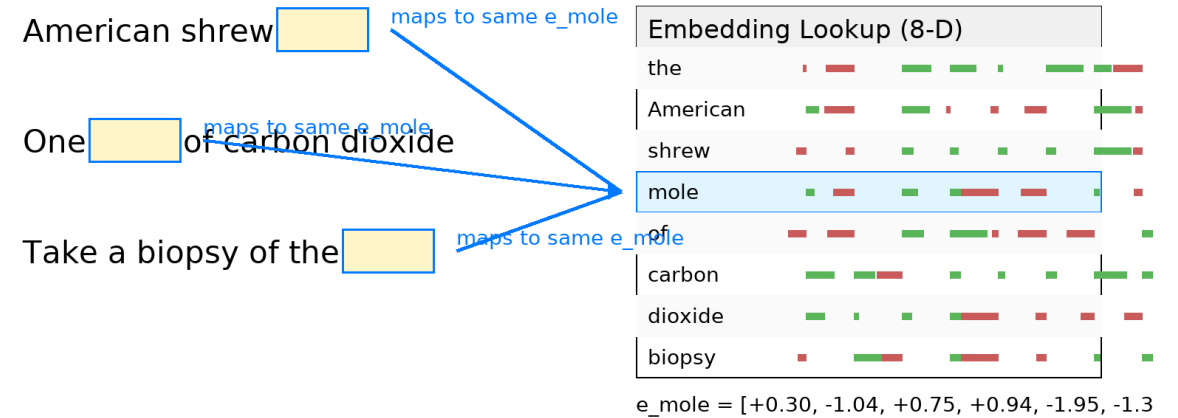# Attention

- By Polisetty Sai Teja

# Recap

- **Embeddings first:** Transformers turn text into high-dimensional vectors called *embeddings*.

- **Lookup, not meaning (yet):** At this stage, an embedding is a learned lookup — the token **"mole"** maps to the same vector in every sentence.

- **What transformers add:** Through **attention**, the model reads surrounding words and **steers** each token's vector toward the meaning required by its context.

- **Why it matters:** Context-aware vectors help the model both **understand neighbors** and **predict the next token** more accurately.

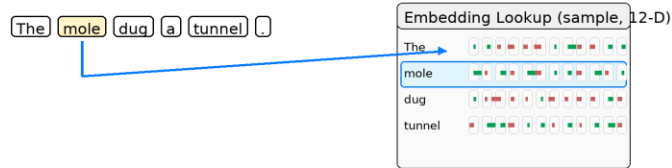**ame Token → Same Initial Embedding (Before Attention**

American shrew [ ]  maps to same e_mole

One [ ] of carbon dioxide  maps to same e_mole

Take a biopsy of the [ ]  maps to same e_mole

Embedding Lookup (8-D)
- the
- American
- shrew
- mole
- of
- carbon
- dioxide
- biopsy

e_mole = [+0.30, -1.04, +0.75, +0.94, -1.95, -1.3

Before attention: every 'mole' maps to the same embedding row (lookup).

# Core Idea of Attention

**Context-free Embeddings (Before Attention)**
Text → tokens → ID → lookup vector (no context yet)

The mole dug a tunnel .

Embedding Lookup (sample, 12-D)
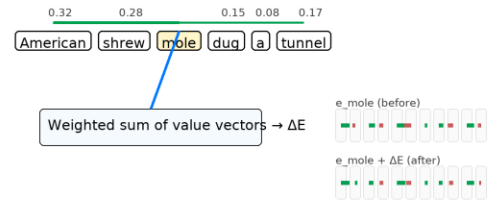
| The |
| mole |
| dug |
| tunnel |

All 'mole' tokens map to the same initial embedding row.

Start point: Tokens begin as context-free embeddings (look up vectors).
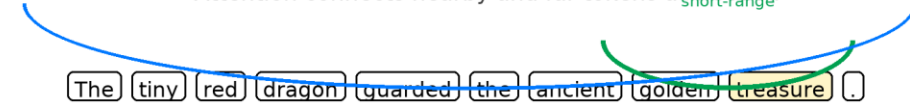
**What Attention Does: Weighted Information Pull**
Each token gathers from others with learned weights → ΔE → new embedding

0.32    0.28        0.15  0.08    0.17
American shrew mole dug a tunnel

Weighted sum of value vectors → ΔE

e_mole (before)

e_mole + ΔE (after)

What attention does: Each token computes how much to pay attention to every other token and pulls in their information as a weighted sum—producing a context-specific update (ΔE) to its embedding.

long-range
**Short- and Long-Range Information Flow**
Attention connects nearby and far tokens as needed
short-range

The tiny red dragon guarded the ancient golden treasure .

Short & long range: These weights connect both nearby words and far-apart words, so information can jump across the sequence(short- and long-distance dependencies).
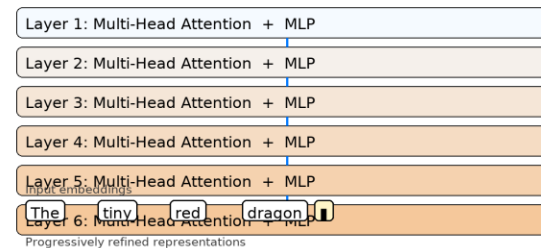
# Core Idea of Attention

**ng Emerges from Context (Same Token, Different S**

American shrew mole ——— 'mole' → animal sense

One mole of carbon dioxide ——— 'mole' → chemistry unit

Take a biopsy of the mole ——— 'mole' → skin lesion

## Next-Token Prediction Uses the Last Position

Final hidden state at last position → Linear → Logits → Softmax → Next token

Input tokens (position indices increasing →)

The | tiny | red | dragon

Final hidden state h.

Linear (Output)
$W\_out \cdot h_t + b \to logits$

Logits (scores)
breathes : 4.2
guards : 3.7
flies : 3.1
sleeps : 2.9
eats : 2.6

$y = softmax(W\_out \cdot h_t + b)$

**Meaning via context:** The same token (e.g., *mole*) is steered to  different regions of the space depending on its neighbours—animal vs chemistry unit vs skin lesion.

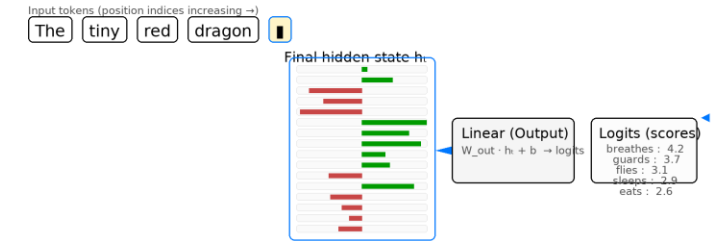## Deep Refinement: [Attention + MLP] × L Layers

Representations grow more abstract with depth

Layer 1: Multi-Head Attention + MLP
Layer 2: Multi-Head Attention + MLP
Layer 3: Multi-Head Attention + MLP
Layer 4: Multi-Head Attention + MLP
Layer 5: Multi-Head Attention + MLP

Input embeddings

The | tiny | red | dragon

Layer 6: Multi-Head Attention + MLP

Progressively refined representations

- **Prediction:** For next-token prediction, the model uses the **final, last-position representation** produced after all layers.

**Deep refinement:** Repeat in attention (plus MLPs) over many layers progressively sharpens meaning.

# Single-Head Attention - What & Why

**Goal:** Let a token "borrow" the right info from other tokens.

**Projections:** $Q_i = E_i W_Q, \ K_i = E_i W_K, \ V_i = E_i W_V$

**Scores:** $s_{ij} = \frac{Q_i \cdot K_j}{\sqrt{d_k}}$ (apply causal **mask** if predicting left→right)

**Weights:** $a_{ij} = \text{softmax}_j(s_{ij})$

**Update:** $\Delta E_i = \sum_j a_{ij} V_j; \ \text{output} = E_i + \Delta E_i$

**Intuition:** $Q$ = what I'm looking for, $K$ = what I offer, $V$ = what to send.

The | hungry | brown | fox | jumped | over | the | sleepy | dog | .

E1 | E2 | E3 | E4 | E5 | E6 | E7 | E8 | E9 | E10

Values: V = E · W_V

| | E1 | E2 | E3 | E4 | E5 | E6 | E7 | E8 | E9 | E10 |
|---|---|---|---|---|---|---|---|---|---|---|
| v1 | 1.00 v1 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| v2 | 0.00 | 1.00 v2 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| v3 | 0.00 | 0.00 | 1.00 v3 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| v4 | 0.00 | 0.50 v2 | 0.40 v3 | 0.10 v4 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| v5 | 0.00 | 0.00 | 0.00 | 0.60 v4 | 0.40 v5 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| v6 | 0.00 | 0.00 | 0.00 | 0.00 | 0.50 v5 | 0.50 v6 | 0.00 | 0.00 | 0.00 | 0.00 |
| v7 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 1.00 v7 | 0.00 | 0.00 | 0.00 | 0.00 |
| v8 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 1.00 v8 | 0.00 | 0.00 | 0.00 |
| v9 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.22 v7 | 0.62 v8 | 0.16 v9 | 0.00 | 0.00 |
| v10 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 1.00 v10 | |

Context contributions (ΔEᵢ) at each position

ΔE1 | ΔE2 | ΔE3 | ΔE4 | ΔE5 | ΔE6 | ΔE7 | ΔE8 | ΔE9 | ΔE10

Add updates to embeddings

E1 + ΔE1 = E1′

E2 + ΔE2 = E2′

E3 + ΔE3 = E3′

E4 + ΔE4 = E4′

E5 + ΔE5 = E5′

E6 + ΔE6 = E6′

E7 + ΔE7 = E7′

E8 + ΔE8 = E8′

E9 + ΔE9 = E9′

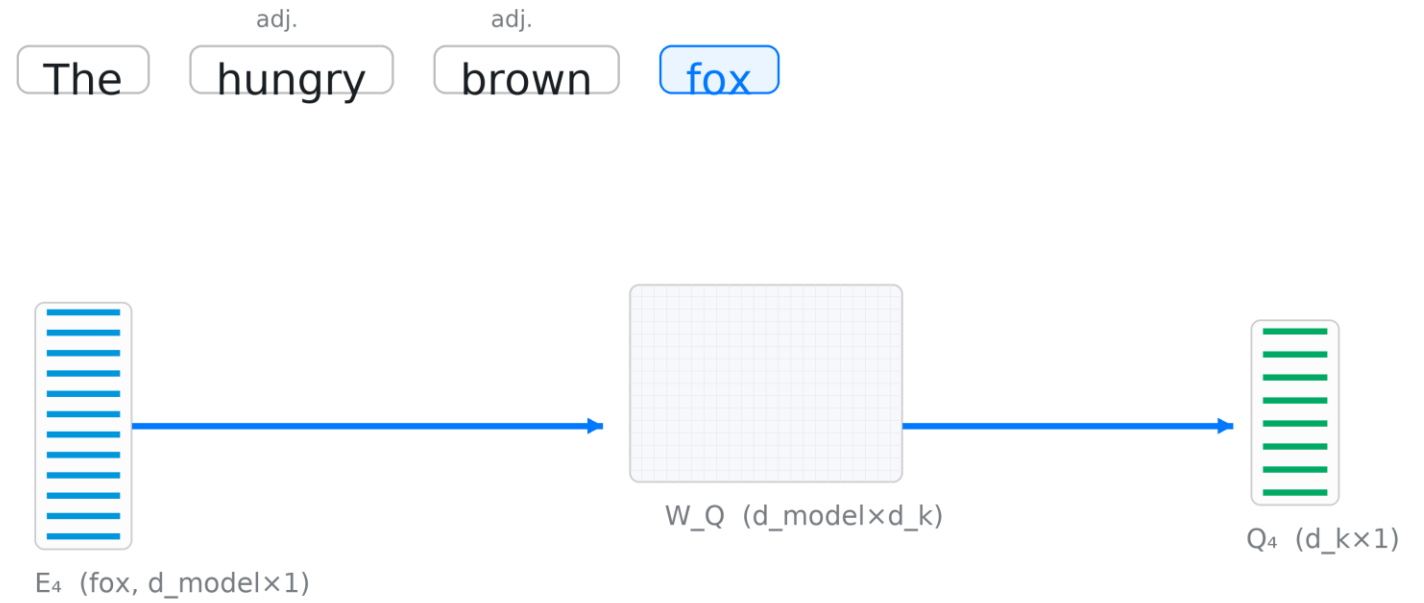E10 + ΔE10 = E10′

- By Sai Teja Polisetty

# Single-Head Attention - Cont...

- **Step 1** — Build the Query (Q)
  - Target token: "fox" in "The hungry brown fox jumped over the sleepy dog."
- **Step 2:** Keys ( $K_i$ ) + "I'm an adjective" tagging
- **Step 3:** Scores(Dot product Q·K → relevance score) + scaling (attention weights - probability distribution) + causal mask(prevents future words from leaking information) → weights
- **Step 4:** Weighted sum of Values → ΔE4 → updated 'fox' - **attention pattern grid**
- **Step 5:** Update the embedding.



The    hungry    brown    fox

① ① Query: Q = E · W_Q

E(fox) → Q(fox)

② ② Keys & Scores: K = E · W_K and s = Q·Kᵀ / √d → a = softmax(s) (masked)

K(The)    K(hungry)    K(brown)

Scores for 'fox' → [The, hungry, brown, fox ...] weights (after softmax)

The   hungry brown fox   ...

③ ③ Values & Update: V = E · W_V and ΔE = Σ a·V → E + ΔE

V(The)
V(hungry)
V(brown)
V(fox)

ΔE → E(fox) + ΔE → context-aware 'fox'

# S1: Query Vector

- Idea: At "fox", the model forms a question: "Are there adjectives before me?"

- **Operation:** $Q_4 = E_4 \cdot W_Q$ projects the "fox" embedding into a **query space** ($d_k$ e.g., 128)

- **Why:** $Q_4$ encodes what "fox" should look for (modifiers like **hungry**, **brown**).

- **Scope:** We compute a **query Qi** for **every** token; we're highlighting Q4 only for clarity.

adj.  adj.

The | hungry | brown | fox

$E_4$ (fox, d_model×1)

$W_Q$ (d_model×d_k)

$Q_4$ (d_k×1)

Build the query for the focus token: $Q_4 = E_4 \cdot W_Q$

# Query vector - cont....

The    hungry    brown    fox    jumped    over    the    sleepy    dog    .

Embeddings  $E_i$  (one per token)

apply to every $E_i$

W_Q  (shared)

$Q_i = E_i \cdot W\_Q$   for i = 1...n

example focus

Q1    Q2    Q3    Q4    Q5    Q6    Q7    Q8    Q9    Q10

Shared projection W_Q maps every embedding $E_i$ into the query space → $Q_i$ (smaller d_k).

# Key vector

**Idea:** Each token advertises what it is. For our head, adjectives like hungry, brown, sleepy emit a "I'm an adjective" signal so nouns (e.g., fox, dog) can find them.
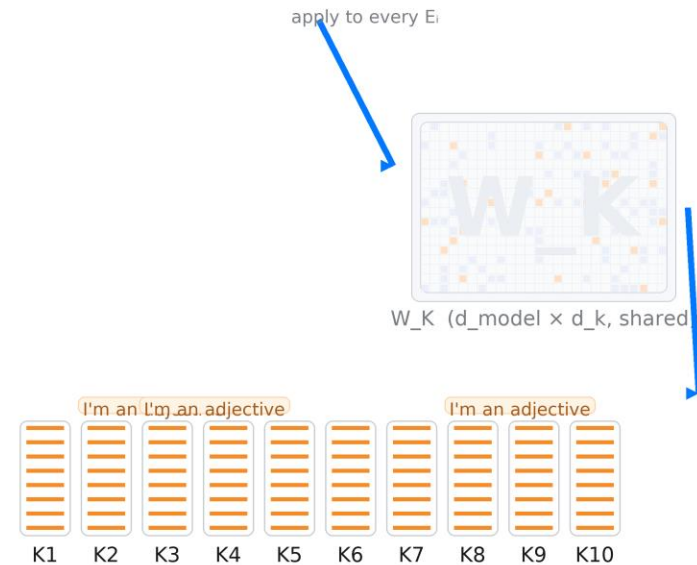
**Operation:** $K_i = E_i \cdot W_k$
- $E_i$ : embedding of token I (size d model)
- $W_k$: learned, shared projection matrix ($d_{model} * d_k$)
- $K_i$: key in the same dk space as queries, so dot products Q.K are comparable

**Why:** Queries need something to match against. Keys are the comparators—learned directions that make the **"right" tokens align with the current token's intent** (e.g., a noun's query aligns with adjective keys). Without keys, the model couldn't score who to attend to.

**Scope**
- We compute a key Ki for every token using the same WK
- Different heads learn different "advertisements" (adjective-ness, subject-verb links, coreference, etc.). This step is independent of masking; masking affects scoring later, not key creation.

The | hungry (adj.) | brown (adj.) | fox | jumped | over | the | sleepy (adj.) | dog | .

Embeddings $E_i$ (one per token)

apply to every $E_i$

W_K (d_model × d_k, shared)

I'm an adjective · I'm an adjective · I'm an adjective

K1 K2 K3 K4 K5 K6 K7 K8 K9 K10

Compute keys for every token: $K_i = E_i \cdot W\_K$ — keys advertise what each token offers to others.

# S2: Key vector

Conceptually, Key vector is **answer** of query vector when both are in the same direction.

# S3: Scores -> Mask -> Weights (queries vs keys)

- **Idea**: Turn "what I'm looking for" (Q) and "what I offer" (K) into match scores, then mask future positions, and finally apply softmax to get attention weights.

- **Operation**:
  - i is the current position (here we illustrate i=4 the word "fox").
  - Mask all positions j>i (future tokens) before softmax in a decoder block.

$$s_j = \frac{Q_i \cdot K_j}{\sqrt{d_k}}$$

$$a_j = \text{softmax}(s)_j$$

| | The Q1 | hungry Q2 | brown Q3 | fox Q4 | jumped Q5 | over Q6 | the Q7 | sleepy Q8 | dog Q9 | . Q10 |
|---|---|---|---|---|---|---|---|---|---|---|
| The K1 | 0.00 | 0.00 | 0.00 | 0.30 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| hungry K2 | 0.00 | 0.00 | 0.00 | 1.20 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| brown K3 | 0.00 | 0.00 | 0.00 | 1.05 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| fox K4 | 0.00 | 0.00 | 0.00 | 0.40 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| jumped K5 | 0.00 | 0.00 | 0.00 | -0.10 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| over K6 | 0.00 | 0.00 | 0.00 | -0.15 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| the K7 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| sleepy K8 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| dog K9 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| . K10 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |

softmax → weights for Q4 ('fox')

| | |
|---|---|
| The | 0.15 |
| hungry | 0.37 |
| brown | 0.32 |
| fox | 0.17 |

Weights sum to 1.00 over visible rows (masked rows e

The | hungry (adj.) | brown (adj.) | fox | jumped | over | the | sleepy (adj.) | dog | .

Target query: $Q_4$ (fox)

Raw scores: $s_j = Q_4 \cdot K_j$

| The | 0.30 |
| hungry | 1.20 |
| brown | 1.05 |
| fox | 0.40 |
| jumped | -0.10 |
| over | -0.15 |
| the | 0.00 |
| sleepy | 0.00 |
| dog | 0.00 |
| . | 0.00 |

Causal mask: positions to the right of 'fox' are ignored (set to $-\infty$ before softmax).

Scale & normalize

$\hat{s}_j = s_j / \sqrt{d_k}$      $a_j = \text{softmax}(\hat{s}_j)$

| The | 0.14 |
| hungry | 0.42 |
| brown | 0.36 |
| fox | 0.08 |

Weights sum to 1.00 over visible positions.

Step 3: scores → mask → softmax. Example for $Q_4$ ('fox'): focus lands on adjectives 'hungry' and 'brown'.

**ATTENTION PATTERN**

## Why
- The dot product measures **alignment** between the current token's intent and other tokens' signals.
- Scaling by sqrt($d_k$) stabilizes training.
- Softmax produces a **probability distribution** (weights) over allowed positions that sum to 1.

## Scope
- Do this for **every position** i.
- Masking policy depends on architecture: **causal** (decoder-only) vs **bidirectional** (encoders don't mask).



Dot size ∝ score (illustrative)   med   high

# ** Context Size

- size of attention pattern is same as that of context size.

- Hence context is not scalable, but to work on it different mechanics came into picture like
  - Sparse attention mechanism
  - Blockwise attention
  - Linformer
  - Reformer
  - Ring attention
  - Long former
  - Adaptive attention span

# S4 : Values -> ΔE -> Updated Embedding

- Idea: Use the attention **weights** to mix **value vectors** $V_j$ . The result is a **context vector** $C_i$ that captures what the current token needs from others. Project it (optional $W_O$ ) and **add** to the original embedding.

- Operation:

$$V_j = E_j W_V$$

$$C_i = \sum_j a_{ij} V_j$$

$$\Delta E_i = C_i W_O \quad \text{(optional head output proj)}$$

$$E'_i = E_i + \Delta E_i \quad \text{(then LayerNorm in the block)}$$

## Why

Values carry the information to transfer (e.g., the adjectives' descriptors). The weights decide how much of each value flows to the current position. Adding ΔEi steers the embedding in a context-aware direction.

## Scope

This happens for every token and for every head in the layer. Multi-head outputs are concatenated and projected before the residual add & LayerNorm.

**Values** $V_j = E_j \cdot W\_V$ **(contributors only)**

V1

The
V2

hungry
V3

brown
V4

fox

**Weights for $Q_4$ (fox)**

| The | 0.14 |
| hungry | 0.42 |
| brown | 0.36 |
| fox | 0.08 |

**Weighted sum → context → update**

$C_4 = 0.14 \cdot v_1 + 0.42 \cdot v_2 + 0.36 \cdot v_3 + 0.08 \cdot v_4$

$\times$ W      $+$      $=$

context $C_4$   $\Delta E_4$      $E_4$      $E_4'$

- By Sai Teja Polisetty
Step 4 (clarified): sum the value vectors with the attention weights to form $C_4$, project to $\Delta E_4$, and add to $E_4 \rightarrow E_4'$.

| | The<br>E1<br>Q1 | adj.<br>hungry<br>E2<br>Q2 | adj.<br>brown<br>E3<br>Q3 | fox<br>E4<br>Q4 | jumped<br>E5<br>Q5 | over<br>E6<br>Q6 | the<br>E7<br>Q7 | adj.<br>sleepy<br>E8<br>Q8 | dog<br>E9<br>Q9 | .<br>E10<br>Q10 |
|---|---|---|---|---|---|---|---|---|---|---|
| E1 The W_V→ v1 | | | | 0.14 v1 | | | | | | |
| E2 hungry W_V→ v2 | | | | 0.42 v2 | | | | | | |
| E3 brown W_V→ v3 | | | | 0.36 v3 | | | | | | |
| E4 fox W_V→ v4 | | | | 0.08 v4 | | | | | | |
| E5 jumped W_V→ v5 | | | | 0.00 v5 | | | | | | |
| E6 over W_V→ v6 | | | | 0.00 v6 | | | | | | |
| E7 the W_V→ v7 | | | | 0.00 v7 | | | | | | |
| E8 sleepy W_V→ v8 | | | | 0.00 v8 | | | | | | |
| E9 dog W_V→ v9 | | | | 0.00 v9 | | | | | | |
| E10 . W_V→ v10 | | | | 0.00 v10 | | | | | | |

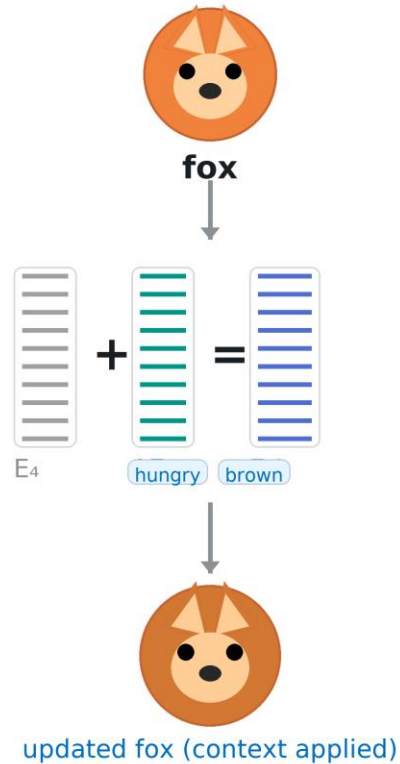Value matrix (white): each row forms $v_i = E_i \cdot W\_V$ (left). The highlighted 'fox' column shows its attention weights × value labels.

- By Sai Teja Polisetty

When we multiply **VALUE MATRIX** with the embedding, we can think of it as saying, if this word is relevant to adjusting the meaning of something else, what exactly should be added to embedding of something else.

**Keys and queries are not needed once u get the attention pattern.

Values: $V_j = E_j \cdot W\_V$
v1 (The)
v2 (hungry)
v3 (brown)
v4 (fox)
v5 (jumped)
v6 (over)
v7 (the)
v8 (sleepy)
v9 (dog)
v10 (.)

Weights for $Q_4$ ('fox')

| The | 0.14 |
| hungry | 0.42 |
| brown | 0.36 |
| fox | 0.08 |
| jumped | 0.00 |
| over | 0.00 |
| the | 0.00 |
| sleepy | 0.00 |
| dog | 0.00 |
| . | 0.00 |

Weighted sum → $\Delta E_4$ → $E_4'$
Contributions $a_j \cdot v_j$ (only non-zero shown)
0.42 × $v_2$ (hungry)
0.36 × $v_3$ (brown)
0.14 × $v_1$ (The)
0.08 × $v_4$ (fox)

sum → context $C_4$

× W_    +    =

$\Delta E_4$    $E_4$    $E_4'$

Step 4: Build V = E·W_V, take the weighted sum using attention weights for $Q_4$ to get context $C_4$, project to $\Delta E_4$ (via W_O), and update $E_4$ → $E_4'$.

E4 = E4 + ΔE4



fox

↓

$E_4$ + (hungry)(brown) = updated fox (context applied)

Attention adds ΔE4 (from hungry/brown), giving E4′ — a co

| | The E1 | hungry E2 | brown E3 | fox E4 | jumped E5 | over E6 | the E7 | sleepy E8 | dog E9 | . E10 |
|---|---|---|---|---|---|---|---|---|---|---|
| The | W_V 0.00 v1 | 0.00 v1 | 0.00 v1 | 0.14 v1 | 0.00 v1 | 0.00 v1 | 0.00 v1 | 0.00 v1 | 0.06 v1 | 0.00 v1 |
| hungry | W_ | 0.00 v2 | 0.00 v2 | 0.42 v2 | 0.00 v2 | 0.00 v2 | 0.00 v2 | 0.00 v2 | | 0.00 v2 |
| brown | W_ | | 0.00 v3 | 0.36 v3 | 0.00 v3 | 0.00 v3 | 0.00 v3 | 0.00 v3 | | 0.00 v3 |
| fox | W_ | | | 0.08 v4 | 0.00 v4 | 0.00 v4 | 0.00 v4 | 0.00 v4 | | 0.00 v4 |
| jumped | W_ | | | | 0.00 v5 | 0.00 v5 | 0.00 v5 | 0.00 v5 | | 0.00 v5 |
| over | W_ | | | | | 0.00 v6 | 0.00 v6 | 0.00 v6 | | 0.00 v6 |
| the | W_ | | | | | | 0.00 v7 | 0.00 v7 | 0.24 v7 | 0.00 v7 |
| sleepy | W_ | | | | | | | 0.00 v8 | 0.60 v8 | 0.00 v8 |
| dog | W_ | | | | | | | | 0.10 v9 | 0.00 v9 |
| . | W_ | | | | | | | | | 0.00 v10 |

ΔE1  ΔE2  ΔE3  ΔE4  ΔE5  ΔE6  ΔE7  ΔE8  ΔE9  ΔE10

Apply updates at every position

E1 +ΔE1 =E1E2 +ΔE2 =E2E3 +ΔE3 =E3E4 +ΔE4 =E4E5 +ΔE5 =E5E6 +ΔE6 =E6E7 +ΔE7 =E7E8

E1′        E2′        E3′        E4′        E5′        E6′        E7′        E8′

This complete process is called = **single head of attention.**

# Connect to the "big model" numbers cont ....

- Total weights = 175,181,291,520
- Organized into 27,938 matrices

| Matrix name | Matrix dimension |
| --- | --- |
| Embedding | D_embed * n_vocab<br>12,288 * 50,257 = 617558016 |
| Key | 12,288 * 128 = 1,572,864 |
| Query | 12,288 * 128 = 1,572,864 |
| Value | 12,288 * 12,288 (wrong)<br>128 * 12,288 = 1,572,864 And<br>12,288 * 128 = 1,572,864 |
| Output | |
| Up-projection | |
| Down-projection | |
| UnEmbedding | n_vocab * D_embed<br>50,257 * 12288 = 617558016 |

# Multi-Head Attention - Big Idea

- Idea
  - Single head attention * 10,000 = parallel ops
  - Each head with its own key/query/value matrices, so it can learn a **different relation** (e.g., adjectives →nouns, subject → verb, long-range links).
  - GPT-3 has 96 attention heads inside each block = 96 attention patterns

Operation:

For each head $h$:

$$Q^{(h)} = E\,W_Q^{(h)}, \quad K^{(h)} = E\,W_K^{(h)}, \quad V^{(h)} = E\,W_V^{(h)}$$

$$A^{(h)} = \mathrm{softmax}\left(\frac{Q^{(h)} K^{(h)\top}}{\sqrt{d_k}}\right), \quad C^{(h)} = A^{(h)} V^{(h)}$$
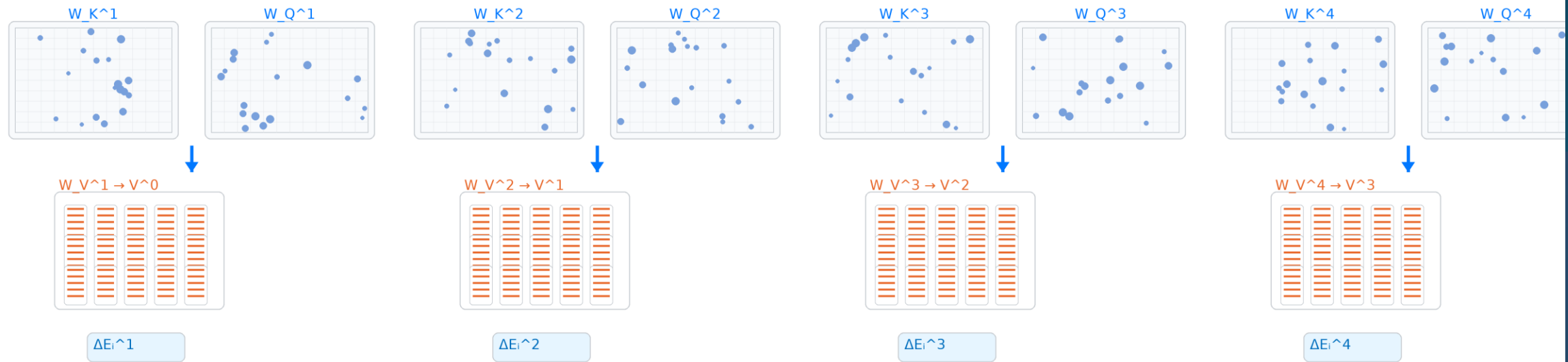
Concat all head outputs and project:

$$\mathrm{MHA}(E) = \left[\, C^{(1)} \,\|\, \cdots \,\|\, C^{(H)} \,\right] W_O$$

Residual add & LayerNorm finish the block:

$$E' = \mathrm{LayerNorm}(E + \mathrm{MHA}(E))$$

**Combine heads: per-head updates add into ΔE$_i$, then E$_i$ ← E$_i$ + ΔE$_i$**

W_K^1          W_Q^1              W_K^2          W_Q^2              W_K^3          W_Q^3              W_K^4          W_Q^4

W_V^1 → V^0                       W_V^2 → V^1                       W_V^3 → V^2                       W_V^4 → V^3

ΔE$_i$^1                          ΔE$_i$^2                          ΔE$_i$^3                          ΔE$_i$^4

Original embedding  E$_i$  +  ΔE$_i$^{(1)}  +  ΔE$_i$^{(2)}  +  ΔE$_i$^{(3)}  +  ΔE$_i$^{(4)}  +  …

→  E$_i$′  (concat heads → W_O, then residual add & LayerNorm)

# Multi-Head Attention – cont...

## Why this works

- Different heads specialize (syntax, agreement, long-distance references, etc.).
- Parallel heads = richer, disentangled signals per token.
- Scales well: GPT-3 uses **96 heads per block** (we illustrate with 4 for clarity).

## Scope

- Heads run **for every token** and **at every layer**.
- The mechanism is identical to single-head; only the number of projections changes.

# Connect to the "big model" numbers cont ....

- Total weights = 175,181,291,520
- Organized into 27,938 matrices

| Matrix name | Matrix dimension |
|---|---|
| Embedding | D_embed * n_vocab<br>12,288 * 50,257 = 617558016 |
| Key | 12,288 * 128 = 1,572,864 * **96** |
| Query | 12,288 * 128 = 1,572,864 * **96** |
| Value | 12,288 * 12,288 (wrong)<br>128 * 12,288 = 1,572,864 * **96** And<br>12,288 * 128 = 1,572,864 * **96** |
| Output | |
| Up-projection | |
| Down-projection | |
| UnEmbedding | n_vocab * D_embed<br>50,257 * 12288 = 617558016 |

- Inside attention, each head uses two "value" projections:
  - **Value↓ ( $W_V$↓)**: projects the model width to the head width
    $W_{V↓} \in R$ (e.g., 12,288→128)
  - After attention mixes values, each head must return to model width via **Value↑**.
    Collecting **all heads' Value↑** blocks into one big projection is the **output matrix**
    $W_O \in R^{(H \cdot dhead) \times dmodel}$.

# Output Matrix

**Output matrix (W_O): the 'Value↑' from ALL heads assembled into one projection**
Docs usually say "value matrix" to mean Value↓ (d_model→d_head).  W_O is the big map back to d_model.

**Value↓**

d_model (e.g., 12,288)    d_head (e.g., 128)

W_V↓ : d_model → d_head

**Value↑**

d_head (e.g., 128)    d_model (e.g., 12,288)

W_V↑ : d_head → d_model

**Stack Value↑ across H hea**

Value↑ (head 1)
Value↑ (head 2)
Value↑ (head 3)
Value↑ (head 4)
Value↑ (head 5)
Value↑ (head 6)
Value↑ (head 7)
Value↑ (head 8)

H × d_head rows  ×  d_model cols

**Output matrix  W_O**

concat

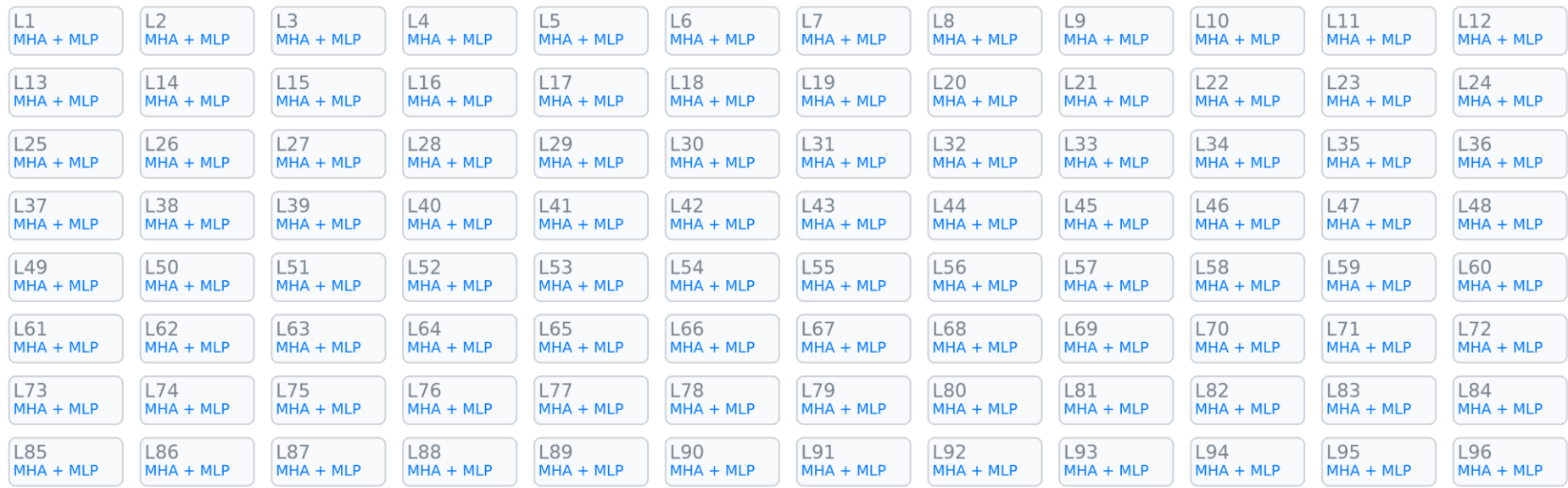shape: (H·d_head) × d_model
Example (GPT-3):  H=96, d_head=12

Takeaway: when someone says "the value matrix," they usually mean Value↓. The 'value up' pieces across heads are assembled into a single Output matrix W_O that maps the concatenated head outpu

- By Sai Teja Polisetty    24

# Beyond attention: the full Transformer block

- Data doesn't stop at attention; each block has Multi-Head Attention → residual+LayerNorm → MLP → residual+LayerNorm.

- The MLP is position-wise (same weights at every token), typically widening by 2–4× then projecting back down.

- Depth builds meaning
  - Blocks are repeated many times (e.g., GPT-3 uses 96 layers).
  - As we go deeper, embeddings absorb richer context from other embeddings that are themselves improving—enabling higher-level abstractions (sentiment, tone, genre), not just grammar.

# Deep Stack: repeated blocks (e.g., 96 layers in GPT-3)

| L1 MHA + MLP | L2 MHA + MLP | L3 MHA + MLP | L4 MHA + MLP | L5 MHA + MLP | L6 MHA + MLP | L7 MHA + MLP | L8 MHA + MLP | L9 MHA + MLP | L10 MHA + MLP | L11 MHA + MLP | L12 MHA + MLP |
| L13 MHA + MLP | L14 MHA + MLP | L15 MHA + MLP | L16 MHA + MLP | L17 MHA + MLP | L18 MHA + MLP | L19 MHA + MLP | L20 MHA + MLP | L21 MHA + MLP | L22 MHA + MLP | L23 MHA + MLP | L24 MHA + MLP |
| L25 MHA + MLP | L26 MHA + MLP | L27 MHA + MLP | L28 MHA + MLP | L29 MHA + MLP | L30 MHA + MLP | L31 MHA + MLP | L32 MHA + MLP | L33 MHA + MLP | L34 MHA + MLP | L35 MHA + MLP | L36 MHA + MLP |
| L37 MHA + MLP | L38 MHA + MLP | L39 MHA + MLP | L40 MHA + MLP | L41 MHA + MLP | L42 MHA + MLP | L43 MHA + MLP | L44 MHA + MLP | L45 MHA + MLP | L46 MHA + MLP | L47 MHA + MLP | L48 MHA + MLP |
| L49 MHA + MLP | L50 MHA + MLP | L51 MHA + MLP | L52 MHA + MLP | L53 MHA + MLP | L54 MHA + MLP | L55 MHA + MLP | L56 MHA + MLP | L57 MHA + MLP | L58 MHA + MLP | L59 MHA + MLP | L60 MHA + MLP |
| L61 MHA + MLP | L62 MHA + MLP | L63 MHA + MLP | L64 MHA + MLP | L65 MHA + MLP | L66 MHA + MLP | L67 MHA + MLP | L68 MHA + MLP | L69 MHA + MLP | L70 MHA + MLP | L71 MHA + MLP | L72 MHA + MLP |
| L73 MHA + MLP | L74 MHA + MLP | L75 MHA + MLP | L76 MHA + MLP | L77 MHA + MLP | L78 MHA + MLP | L79 MHA + MLP | L80 MHA + MLP | L81 MHA + MLP | L82 MHA + MLP | L83 MHA + MLP | L84 MHA + MLP |
| L85 MHA + MLP | L86 MHA + MLP | L87 MHA + MLP | L88 MHA + MLP | L89 MHA + MLP | L90 MHA + MLP | L91 MHA + MLP | L92 MHA + MLP | L93 MHA + MLP | L94 MHA + MLP | L95 MHA + MLP | L96 MHA + MLP |

As depth increases:
- lower: local grammar / short-range links
- middle: long-range dependencies, entities
- higher: topic, sentiment, tone, genre (poem vs. research)

**Transformer Block: Multi-Head Attention + MLP (with residuals & LayerNorm)**
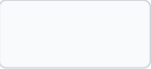
E (token embeddings, all positions in parallel)

Multi-Head Attention

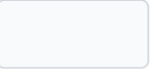Q,K,V per head → attention → concat → W_O

Residual Add + LayerNorm

Position-wise MLP

2–4× width up-proj → nonlinearity → down-proj

Residual Add + LayerNorm

Output of block → feeds next block (same for every token).

Parallelism:
- Across tokens (matrix multiplies over the whole sequence)
- Across heads (dozens of heads per block)
→ GPUs exploit this for scale and throughput

# Connect to the "big model" numbers cont ....

- Total weights = 175,181,291,520

- Organized into 27,938 matrices

| Matrix name | Matrix dimension |
|---|---|
| Embedding | D_embed * n_vocab<br>12,288 * 50,257 = 617558016 |
| Key | 12,288 * 128 = 1,572,864 * 96 *** 96 layers** |
| Query | 12,288 * 128 = 1,572,864 * 96 *** 96 layers** |
| Value | 12,288 * 12,288 (wrong)<br>128 * 12,288 = 1,572,864 * 96 *** 96 layers** And<br>12,288 * 128 = 1,572,864 * 96 *** 96 layers** |
| Output | |
| Up-projection | |
| Down-projection | |
| UnEmbedding | n_vocab * D_embed<br>50,257 * 12288 = 617558016 |

**Total studied params** so far: ~58B parameters

**Yet to study params:** ~ 117B parameters

# What's next? - MLP layers ...

**Where the parameters live (attention vs other)**

Attention params (K/Q/V + Output)
≈ 57,982,058,496  (≈58B)  ≈ one-third

Other params (MLPs, embeddings, layer norms, etc.)
≈ two-thirds of total

Takeaways

• Attention is powerful for routing information, and it parallelizes well across tokens & heads.

• Most parameters come from the MLPs (big up/down projections at each layer).

• Scaling works because GPUs handle the large matrix multiplies in parallel.