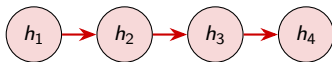# The Transformer Architecture

Self-Attention · Multi-Head Attention · Positional Encoding · Encoder–Decoder
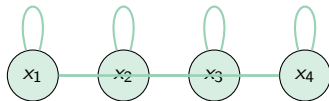
# The key idea

**LSTM**



Sequential: $O(n)$ path

**Transformer**



Parallel: $O(1)$ path

**"Attention Is All You Need"** (Vaswani et al., 2017)

**Parallel**
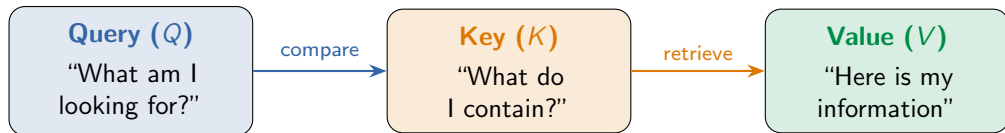Process all tokens simultaneously

**Direct access**
Any token can attend to any other

**Contextual**
Representations change with context

# Attention — intuition

**Think of it as a soft dictionary lookup**

| Query ($Q$) | | Key ($K$) | | Value ($V$) |
|---|---|---|---|---|
| "What am I looking for?" | compare → | "What do I contain?" | retrieve → | "Here is my information" |

Example: translating "Le chat dort" → "The cat sleeps"

"sleeps"    0.05    0.10    0.85  "Le"    "chat"    "dort"

Query                          Keys

Output = weighted sum of values: $0.05 \cdot v_{\text{Le}} + 0.10 \cdot v_{\text{chat}} + 0.85 \cdot v_{\text{dort}}$

# Scaled dot-product attention

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

$QK^\top$
Similarity scores

$\div \sqrt{d_k}$
Scale down

softmax
Normalize

$\times V$
Weighted sum

**Why $\sqrt{d_k}$?**
Without scaling, dot products grow with $d_k$, pushing softmax into saturation (near 0 or 1) where gradients vanish

**Dimensions:**
$Q$: ($n \times d_k$)    $K$: ($m \times d_k$)
$V$: ($m \times d_v$)
$QK^\top$: ($n \times m$)    Output: ($n \times d_v$)

Each output row is a **weighted average** of value vectors, where weights come from query-key similarity

# Attention — worked example

**Tokens:** The cat sat

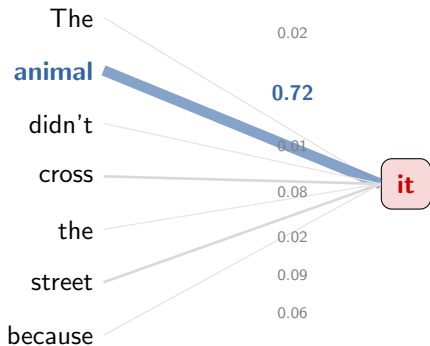|  | $QK^\top / \sqrt{d_k}$ | | | Weights | | |
|---|---|---|---|---|---|---|
|  | The | cat | sat | | | |
| The | 1.2 | 0.5 | 0.3 | .49 | .28 | .23 |
| cat | 0.4 | 2.1 | 0.8 | .12 | .63 | .25 |
| sat | 0.2 | 0.9 | 1.5 | .15 | .30 | .55 |

softmax →

Row for "cat": attends 63% to itself, 25% to "sat", 12% to "The"

Output for "cat" $=$
$0.12 \cdot v_{\text{The}} + 0.63 \cdot v_{\text{cat}} + 0.25 \cdot v_{\text{sat}}$

Each output is a **context-aware** representation — unlike Word2Vec, the same word gets different embeddings in different contexts

# Attention learns meaningful relationships

"The **animal** didn't cross the street because **it** was too tired"



The
**animal**
didn't
cross
the
street
because

0.02
**0.72**
0.01
0.08
0.02
0.09
0.06

**it**

**Coreference resolved!**

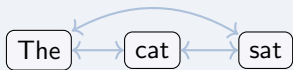"it" attends most strongly to "animal" — the model learned that "it" refers to "animal"

**No explicit rule**

This emerges from training — the model learns which tokens are relevant for each query

Different attention heads specialize in different relationships (see next slides)

# Self-attention vs. cross-attention



**Self-Attention**
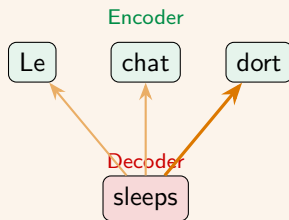
The ↔ cat ↔ sat

$Q$, $K$, $V$ all come from
the **same** sequence

Used in: encoder, decoder (masked)

**Cross-Attention**

Encoder

Le    chat    dort

Decoder
sleeps

$Q$ from decoder
$K$, $V$ from encoder

Used in: decoder (enc-dec models)

# Multi-head attention



```
Input X
```

| Head 1 | Head 2 | Head 3 | ... | Head h |
| $W_Q^1, W_K^1, W_V^1$ | $W_Q^2, W_K^2, W_V^2$ | $W_Q^3, W_K^3, W_V^3$ | | $W_Q^h, W_K^h, W_V^h$ |

Concatenate

Linear $W^O$

Output

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \ldots, \text{head}_h) \, W^O$$
$$\text{where head}_i = \text{Attention}(QW_Q^i, KW_K^i, VW_V^i)$$

# What different heads learn

Different heads specialize in different types of relationships:
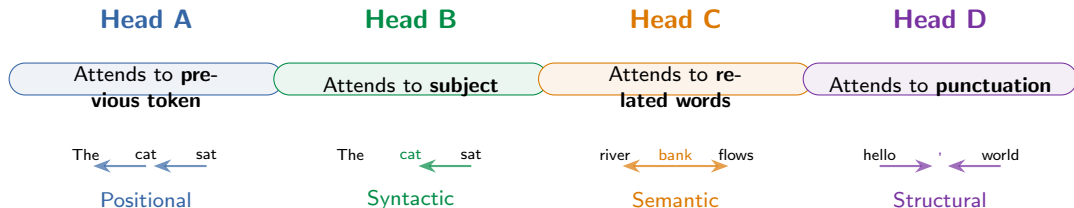
| Head A | Head B | Head C | Head D |
|--------|--------|--------|--------|
| Attends to **previous token** | Attends to **subject** | Attends to **related words** | Attends to **punctuation** |

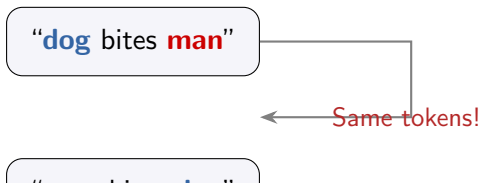| The  cat  sat | The  cat  sat | river  bank  flows | hello  '  world |
|---------------|---------------|--------------------|-----------------|
| Positional | Syntactic | Semantic | Structural |

**Multiple heads** = multiple "perspectives" on the same input.

Typical: $h = 8$ (BERT-base) or $h = 12$–$96$ (larger models). $\quad d_k = d_{model}/h$

Total compute is the same as single-head attention with full $d_{model}$, since each head uses $d_k = d_{model}/h$

# Positional encoding — why we need it

"**dog** bites **man**"

~~Same tokens!~~

"~~man bites dog~~"

**Problem:** Self-attention computes
dot products between token pairs.

It's **permutation-invariant** — swapping
token order doesn't change the output!

But word order matters: these two sen-
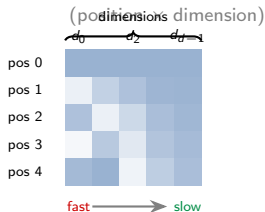tences mean very different things.

**Solution:** Add **positional information** to the input embeddings

$$\text{input}_i = \text{token\_embedding}_i + \text{position\_encoding}_i$$

# Sinusoidal positional encoding

$$PE(pos, 2i) = \sin\left(\frac{pos}{10000^{2i/d}}\right)$$
$$PE(pos, 2i + 1) = \cos\left(\frac{pos}{10000^{2i/d}}\right)$$

**PE matrix**

(positions dimension)

$d_0$ $d_2$ $d_{d-1}$

pos 0
pos 1
pos 2
pos 3
pos 4

fast ——→ slow

**Unique pattern**
Each position gets
a distinct encoding

**Relative positions**
$PE_{pos+k}$ can be
expressed as linear
function of $PE_{pos}$

**No learned params**
Fixed, deterministic,
works for any length

**Multi-scale**
Low dims = fine
position, high
dims = coarse

$input_i = embedding(x_i) + PE(i)$   (element-wise addition)

# Modern positional encodings

## Sinusoidal

Vaswani et al., 2017

Added to embeddings
Absolute position
Fixed (not learned)

Used by: original Transformer

## RoPE

Su et al., 2021

Rotates $Q$
and $K$ vectors
Relative position
No extra parameters

Used by: LLaMA, Mistral

## ALiBi

Press et al., 2022

Bias in attention scores
Linear distance penalty
No extra parameters

Used by: BLOOM, MPT

**Where position info is injected:**

Sinusoidal        RoPE        ALiBi

Embedding $\longrightarrow$ $Q, K, V$ $\longrightarrow$ $QK^\top$ $\longrightarrow$ softmax $\longrightarrow$ output

**RoPE** is the most popular today — it naturally encodes relative position and supports length extrapolation with techniques like YaRN

# The Transformer block



$\times$ **N** layers

Output

**Layer normalization**
Stabilizes training
Pre-norm (GPT) vs post-norm (original)

Add & LayerNorm

**Feed-Forward (per token)**
$\text{FFN}(x) = W_2 \cdot \text{GELU}(W_1 x + b_1) + b_2$
Hidden dim: $4 \times d_{\text{model}}$

Feed-Forward Network

residual

**Residual connections**
output $= x + \text{sublayer}(x)$
Enables deep stacking (gradient flow)

Add & LayerNorm

residual

Multi-Head Attention

**Multi-Head Attention**
Captures token relationships (self or cross-attention)

Input $X$

# The Transformer encoder

**Bidirectional**

Every token attends to every other token (no masking)

**Used by:**

BERT, RoBERTa, DeBERTa

Understanding tasks: classification, NER, QA

Contextual representations

**Encoder Stack**

Self-Attention

Add & Norm

Feed-Forward

Add & Norm

$\times N$

Tokens → Embedding → + ← Pos. Enc.

# The Transformer decoder

## Decoder Block

Masked Self-Attention

Add & Norm

Cross-Attention

Add & Norm

Feed-Forward

Add & Norm

$\times N$

From encoder
($K$, $V$)

**Decoder-only** models (GPT) skip cross-attention — only masked self-attention + FFN

## Causal Mask

attend    masked

Token $t$ can only
see tokens $\leq t$

**Used by:**

GPT (decoder-only)
T5 decoder

Generation tasks

# Encoder–Decoder architecture



**Encoder**
- Self-Attention
- Feed-Forward
- × *N*

Bonjour le monde

$K, V$

$P(\text{next token})$

Linear + Softmax

**Decoder**
- Masked Self-Attn
- Cross-Attention
- Feed-Forward
- × *N*

Hello the world

**Used for:** machine translation, summarization, question answering
**Models:** T5, BART, mBART, original Transformer

# Three Transformer architectures

## Encoder-Only

Bidirectional
MLM training
Understanding tasks

BERT, RoBERTa

## Decoder-Only

Autoregressive (L→R)
CLM training
Generation tasks

GPT, LLaMA, Mistral

## Encoder–Decoder

Bidirectional enc. +
autoregressive dec.

Seq-to-seq tasks

T5, BART, mBART

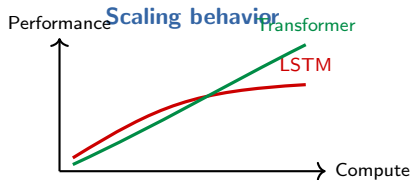|  | Enc-Only | Dec-Only | Enc-Dec |
|---|---|---|---|
| Direction | Bidirectional | Left-to-right | Both |
| Best for | Classification, NER | Text generation | Translation, summary |
| Today's trend | Less common | **Dominant** | Niche |

The trend since GPT-3 (2020): **decoder-only** models at scale
can do almost everything, including understanding tasks.

# Why Transformers won

| | LSTM | Transformer |
| --- | --- | --- |
| Parallelizable | No ($h_t$ needs $h_{t-1}$) | Yes (all tokens at once |
| Long-range path | $O(n)$ steps | $O(1)$ (direct attention) |
| Context | Limited by hidden size | Full context window |
| Scaling | Diminishing returns | Log-linear improvement |
| Training speed | Slow (sequential) | Fast (GPU-friendly) |

**Scaling behavior**



**The scaling insight:**

Transformers reliably improve with
more data, parameters, and compute

This enabled the LLM revolution:
GPT-3, PaLM, LLaMA, Claude, . . .

# Transformer dimensions in practice

| Model | Layers | $d_{model}$ | Heads | $d_{ff}$ | Params |
|-------|--------|-------------|-------|----------|--------|
| BERT-base | 12 | 768 | 12 | 3072 | 110M |
| BERT-large | 24 | 1024 | 16 | 4096 | 340M |
| GPT-2 | 12 | 768 | 12 | 3072 | 117M |
| GPT-3 | 96 | 12288 | 96 | 49152 | 175B |
| LLaMA-2 7B | 32 | 4096 | 32 | 11008 | 7B |
| LLaMA-2 70B | 80 | 8192 | 64 | 28672 | 70B |

$d_k = \frac{d_{model}}{h}$

Head dimension

$d_{ff} \approx 4 \times d_{model}$

FFN hidden size

$Params \approx 12 \cdot N \cdot d^2$

Rough estimate

# Further reading

**Attention & Transformers**

- Vaswani et al. (2017), "Attention Is All You Need" — the original Transformer paper
- Bahdanau et al. (2015), "Neural Machine Translation by Jointly Learning to Align and Translate"
- Jay Alammar, "The Illustrated Transformer"

**Positional Encodings**

- Su et al. (2021), "RoFormer: Enhanced Transformer with Rotary Position Embedding" (RoPE)

**Architecture Variants & Surveys**

- Devlin et al. (2019), "BERT: Pre-training of Deep Bidirectional Transformers" (encoder-only)
- Radford et al. (2018/2019), "Improving/Language Models are Unsupervised Multi-task Learners" (GPT/GPT-2)
- Lin et al. (2022), "A Survey of Transformers" — comprehensive taxonomy

# Questions?

Next: Early Notable Models — GPT, BERT, T5