<u>Attn</u>

query   key$_1$   value$_1$
          k$_2$    V$_2$
          k$_3$    V$_3$
          ⋮        ⋮

$$out = \sum_i sim(q, k_i) V_i$$

$$Out = softmax\left(\frac{QK^T}{\lambda}\right)V$$

Self-attn:

$k^{(2)} q^{(2)} v^{(2)}$



$k_i = W_k S_i$

$q_i = W_q S_i$

$V_i = W_v S_i$

$W_k^{(2)}$   $W_q^{(2)}$ $W_v^{(2)}$ S ← Multi-head

$W_k$ , $W_q$ , $W_v$

S'

Attn

X-attn : q come from a diff. seq.

<u>Transformers</u>   Encoder   Decoder



Target   Eos

LN

FCN

LN

Xattn

LN

Self attn

Mask

q

k

k,V

+ pos encoding
Input seq

+ pos encoding
sos   Target

the dog

BERT     ← TS →     GPT

LN

FCN

LN

self attn

+ pos encoding

TS   DONE

the | dog

BERT

sequential?

S1   S2

CLS   SEP

X

the | dog

Feature
learning

QA,

Generation

Seq   predict
next
token

GPT

Seq

---

Inputs →
(big)

NN

AE encoder

AE
Decoder

**Proxy task**
— reconstruction
— Masked prediction
— contrastive (similarity clustering)
— classification
⋮

Task
data →

NN base

New
pred
head

→ Preds

**What to
Finetune**

Task
data

→ pred

Q₁  A₁  Q₂  A₂  Q₃  (output)
A₃

variable

| Feature Extraction (Linear probing) | Full finetuning (partial OK too) | Hard Prompting | Soft Prompting | Everything Else |
|---|---|---|---|---|
| ← Almost every model ↗ | Best perf | No (or minimal) training data | Doesn't need vp good pretrained models | |
| Moderate amount of data | — Lots of data — one model/task = Memory | Capped performance Mostly for LLM | Good for batches w mixed tasks Needs moderate amount | |

Not the best perf

of data

Doesn't work for all
models
(mostly bb transformers)

New vocab items    Embeddings

| Are | birds | mammals | ? |

Transformer

False