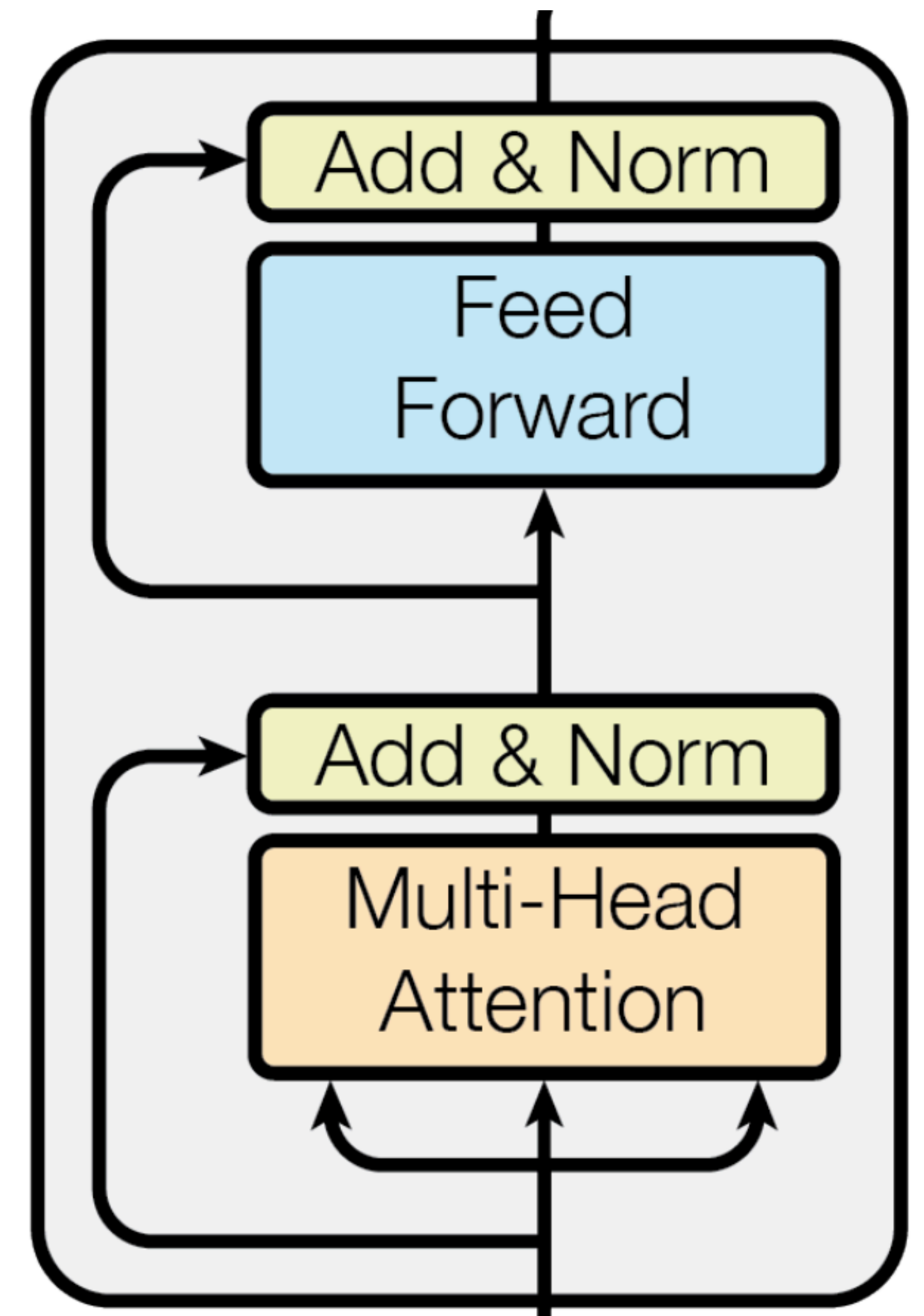# Transformer Architecture
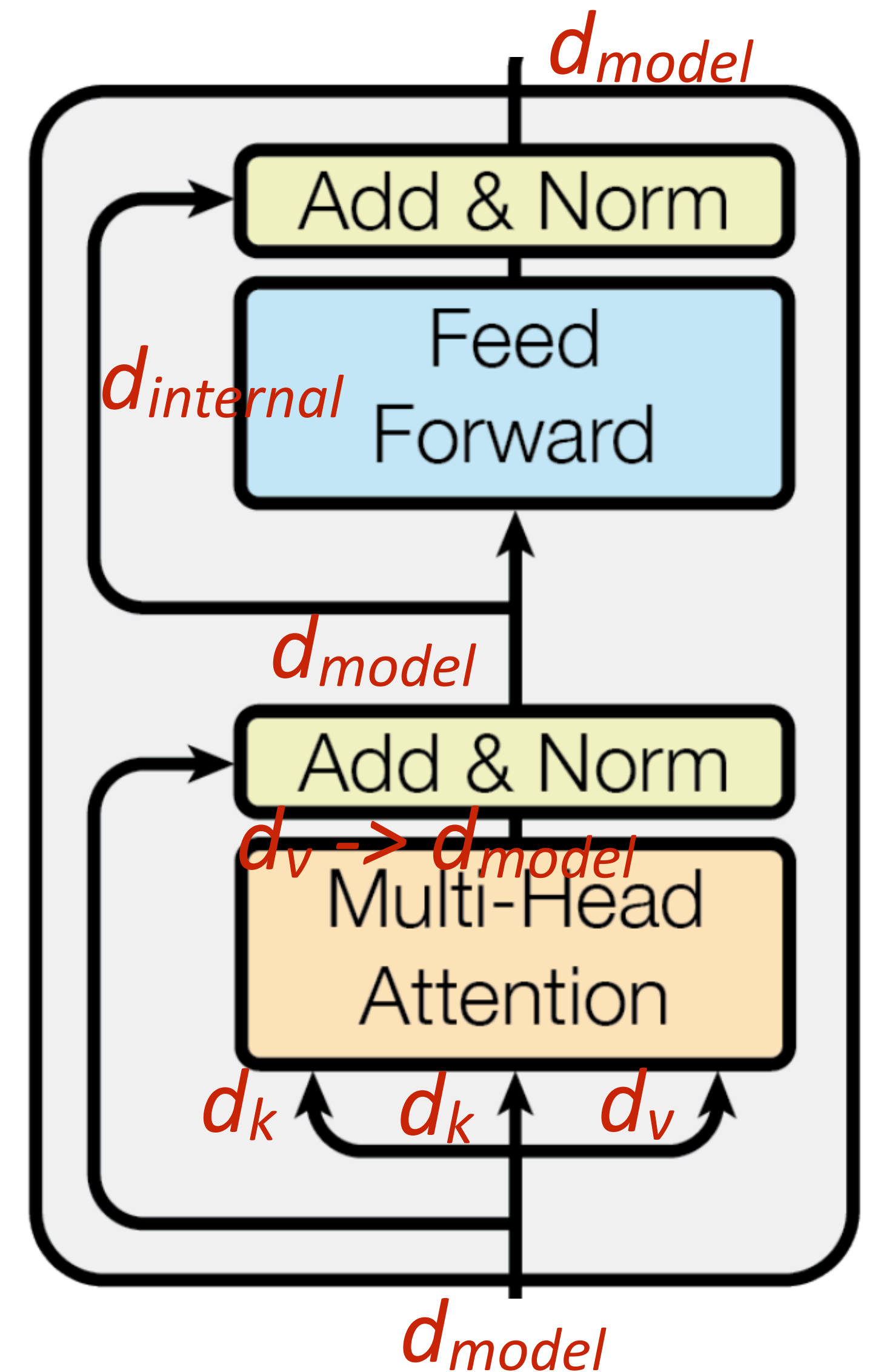
- Self-attention is not the whole story; we'll describe what goes into a Transformer

- Alternate multi-head self-attention with feedforward layers that **operate over each word individually**

$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2$$

  - These feedforward layers are where most of the parameters are

- Residual connections in the model: input of a layer is added to its output

- Layer normalization: controls the scale of different layers in very deep networks

# Transformer Architecture

▸ Vectors: $d_{model}$

▸ Queries/keys: $d_k$, always smaller than $d_{model}$

▸ Values: separate dimension $d_v$, output is multiplied by $W^O$ which is $d_v \times d_{model}$ so we can get back to $d_{model}$ before the residual

▸ FFN can explode the dimension with $W_1$ and collapse it back with $W_2$

$$\mathrm{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2$$
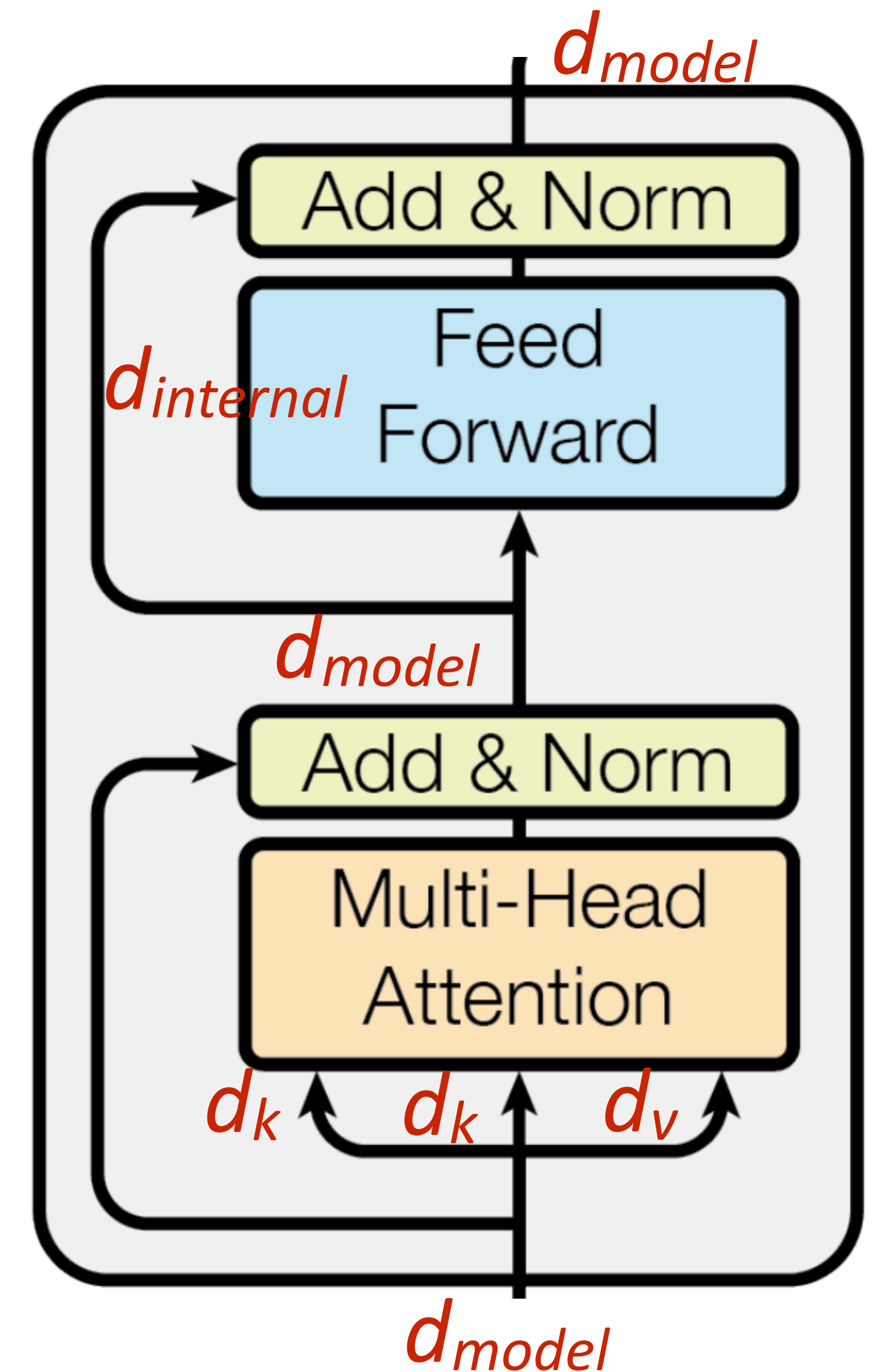


Vaswani et al. (2017)

# Transformer Architecture: Sizes

|  | $N$ | $d_{\text{model}}$ | $d_{\text{ff}}$ | $h$ | $d_k$ | $d_v$ |
|---|---|---|---|---|---|---|
| base | 6 | 512 | 2048 | 8 | 64 | 64 |

‣ From Vaswani et al.

| Model Name | $n_{\text{params}}$ | $n_{\text{layers}}$ | $d_{\text{model}}$ | $n_{\text{heads}}$ | $d_{\text{head}}$ |
|---|---|---|---|---|---|
| GPT-3 Small | 125M | 12 | 768 | 12 | 64 |
| GPT-3 Medium | 350M | 24 | 1024 | 16 | 64 |
| GPT-3 Large | 760M | 24 | 1536 | 16 | 96 |
| GPT-3 XL | 1.3B | 24 | 2048 | 24 | 128 |
| GPT-3 2.7B | 2.7B | 32 | 2560 | 32 | 80 |
| GPT-3 6.7B | 6.7B | 32 | 4096 | 32 | 128 |
| GPT-3 13B | 13.0B | 40 | 5140 | 40 | 128 |
| GPT-3 175B or "GPT-3" | 175.0B | 96 | 12288 | 96 | 128 |

‣ From GPT-3; $d_{head}$ is our $d_k$



Vaswani et al. (2017)

# Transformer Architecture: Sizes

| 1 | description | FLOPs / update | % FLOPS MHA | % FLOPS FFN | % FLOPS attn | % FLOPS logit |
|---|---|---|---|---|---|---|
| 8 | OPT setups | | | | | |
| 9 | 760M | 4.3E+15 | 35% | 44% | 14.8% | 5.8% |
| 10 | 1.3B | 1.3E+16 | 32% | 51% | 12.7% | 5.0% |
| 11 | 2.7B | 2.5E+16 | 29% | 56% | 11.2% | 3.3% |
| 12 | 6.7B | 1.1E+17 | 24% | 65% | 8.1% | 2.4% |
| 13 | 13B | 4.1E+17 | 22% | 69% | 6.9% | 1.6% |
| 14 | 30B | 9.0E+17 | 20% | 74% | 5.3% | 1.0% |
| 15 | 66B | 9.5E+17 | 18% | 77% | 4.3% | 0.6% |
| 16 | 175B | 2.4E+18 | 17% | 80% | 3.3% | 0.3% |

Credit: Stephen Roller on Twitter

# Transformer: Complete Model



‣ Original Transformer paper presents an **encoder-decoder** model

‣ Right now we don't need to think about both of these parts — will return in the context of sequence-to-sequence models later

‣ Can turn the encoder into a decoder-only model through use of a triangular causal attention mask (only allow attention to previous tokens)

Vaswani et al. (2017)