# LLaMA 3
## What are the dimensions?

# 8B model parameters

- # of Layers = __32__

- # of Attention Heads = __32__

- # of Vocabulary Words = __128K__

- # of Feature Dimensions = __4096__

- # of Hidden Dimensions = __5012__

- Context Window Size = __8K__

$$4096 \cdot 1.3 = 5012$$

⌥ main ⌄   Meta-Llama-3-8B / original / params.json ⧉

pcuenq  **HF STAFF**   Upload original checkpoint (#1)

</> raw   ⏱ history   ☺ blame   ✎ contribute   🗑 delete   ⊘ No virus

```
1    {
2        "dim": 4096,        ←
3        "n_layers": 32,     ←
4        "n_heads": 32,      ←
5        "n_kv_heads": 8,
6        "vocab_size": 128256,  ←
7        "multiple_of": 1024,
8        "ffn_dim_multiplier": 1.3,  ←
9        "norm_eps": 1e-05,
10       "rope_theta": 500000.0
11   }
```

AI by Hand ✍ 2024 © Tom Yeh

# Llama

```
35    ∨    class Llama:
36             @staticmethod
37    ∨        def build(
38                 ckpt_dir: str,
39                 tokenizer_path: str,
40                 max_seq_len: int,
41                 max_batch_size: int,
42                 model_parallel_size: Op
43                 seed: int = 1,
44             ) -> "Llama":
```

8K

# Transformer

```
251   v    class Transformer(nn.Module):
252   v        def __init__(self, params:
253               super().__init__()
254               self.params = params
255               self.vocab_size = param
256               self.n_layers = params.
257
```

128 K

32

# Self Attention

```
90   ∨   class Attention(nn.Module):
91   ∨       def __init__(self, args: ModelArgs):
92               super().__init__()
93               self.n_kv_heads = args.n_heads if args.n_k
94               model_parallel_size = fs_init.get_model_pa
95               self.n_local_heads = args.n_heads // model_
96               self.n_local_kv_heads = self.n_kv_heads //
97               self.n_rep = self.n_local_heads // self.n_
98               self.head_dim = args.dim // args.n_heads
```

128

4096

32

AI by Hand ✍️ 2024 © Tom Yeh

# Feed Forward

```
193  ∨    class FeedForward(nn.Module):
194  ∨        def __init__(
195                self,
196                dim: int,
197                hidden_dim: int,
198                multiple_of: int,
199                ffn_dim_multiplier: Optio
200            ):
```
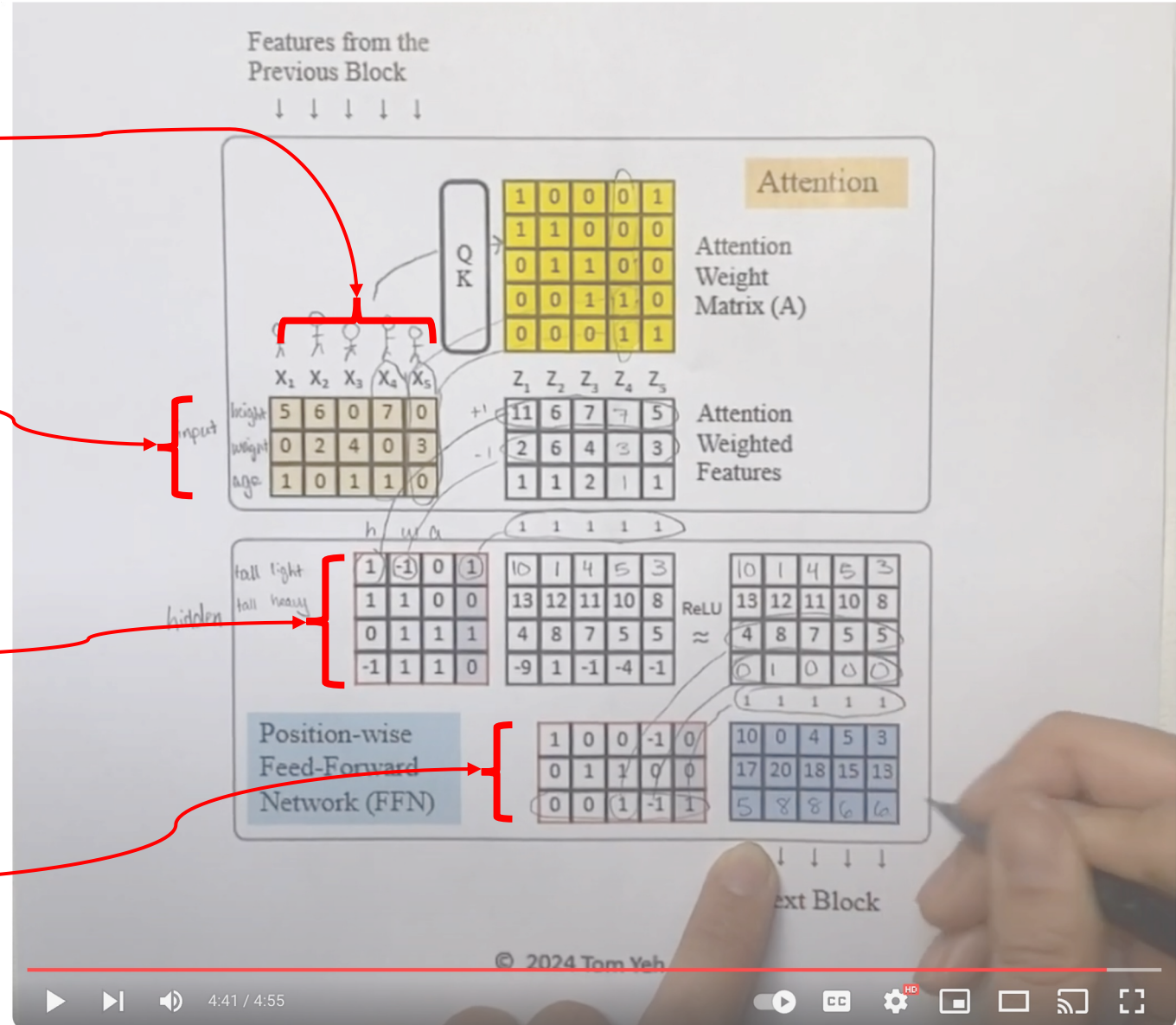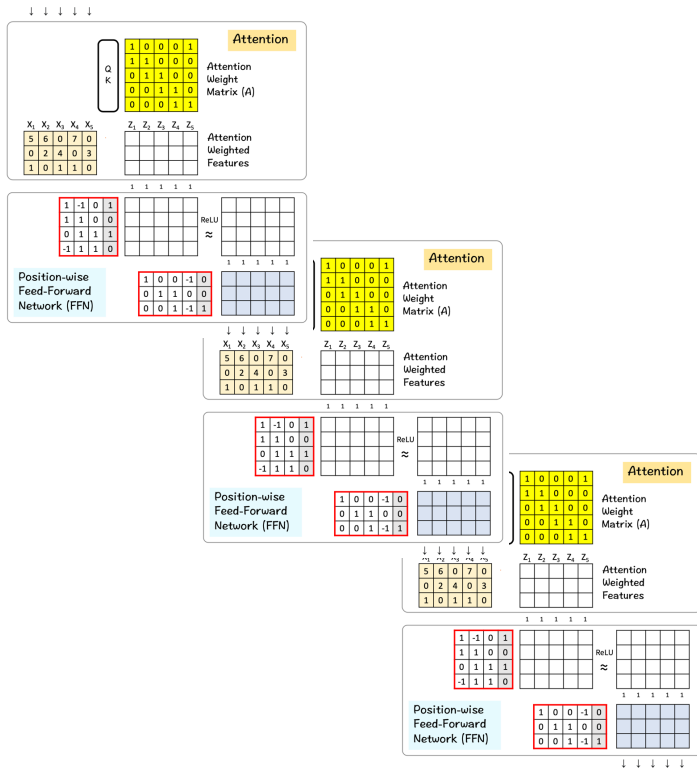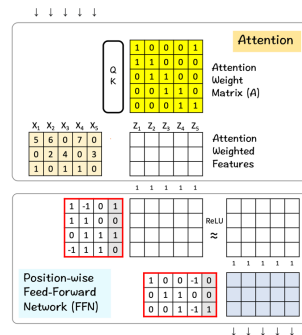
4096

5012

# Transformer Block



8K

4096

5012

4096

**S2.E1 - Transformer - AI by Hand ✍️ with Anna**

# Layers



x  32

AI by Hand ✍️ 2024 © Tom Yeh

# Input / Output



**8K**

**128K**

**4096**

**32**

**4096**

**128K**

Attention

| Q | 1 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| K | 1 | 1 | 0 | 0 | 0 |
|   | 0 | 1 | 1 | 0 | 0 |
|   | 0 | 0 | 1 | 1 | 0 |
|   | 0 | 0 | 0 | 1 | 1 |

Attention Weight Matrix (A)

| $X_1$ | $X_2$ | $X_3$ | $X_4$ | $X_5$ |
|---|---|---|---|---|
| 5 | 6 | 0 | 7 | 0 |
| 0 | 2 | 4 | 0 | 3 |
| 1 | 0 | 1 | 1 | 0 |

$Z_1$ $Z_2$ $Z_3$ $Z_4$ $Z_5$

Attention Weighted Features

1 1 1 1 1

| 1 | -1 | 0 | 1 |
|---|---|---|---|
| 1 | 1 | 0 | 0 |
| 0 | 1 | 1 | 1 |
| -1 | 1 | 1 | 0 |

ReLU ≈

1 1 1 1 1

Position-wise Feed-Forward Network (FFN)

| 1 | 0 | 0 | -1 | 0 |
|---|---|---|---|---|
| 0 | 1 | 1 | 0 | 0 |
| 0 | 0 | 1 | -1 | 1 |

**X**