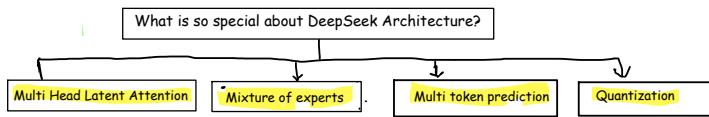
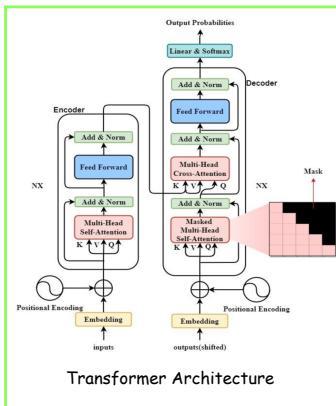


# DeepSeek Architecture



## 1. Multi Head Latent Attention

### 1. Multi Head Attention



### Encoder Block

Input Sentence: "Transformers are amazing"

1. Tokenization: Split the raw sentence into tokens using a tokenizer (e.g., BPE, Sentence piece, word-piece).

Ex.

```
python ⌂ Copy ⌂ Edit
tokens = ["Transformers", "are", "amazing"]
```

2. Vocabulary Lookup (Token → Token ID): Use a predefined vocabulary to map each token to a unique integer ID.

Ex.

```
Suppose the vocabulary has:
python ⌂ Copy ⌂ Edit
"Transformers" -> 17
"are" -> 56
"amazing" -> 94

Result:
python ⌂ Copy ⌂ Edit
token_ids = [17, 56, 94]
```

3. Token Embedding (Learned)

Use a learnable embedding matrix  $E \in \mathbb{R}^{(|V| \times d_{\text{model}})}$   
Each token ID  $t$  maps to a vector  $E[t] \in \mathbb{R}^{d_{\text{model}}}$

```
python ⌂ Copy ⌂ Edit
import torch.nn as nn

embedding_layer = nn.Embedding(num_embeddings=30000, embedding_dim=512)
token_embeddings = embedding_layer(torch.tensor([17, 56, 94])) # shape: (3, 512)

Now:
mathematica ⌂ Copy ⌂ Edit
TokenTable = {
  {17}, # "Transformers"
  {56}, # "are"
  {94} # "amazing"
} -> shape: {seq_len, d_model}w512
```

4. Positional Encoding: Transformers have no recurrence — so we add position info.

Compute fixed sinusoidal encoding:  
For each position  $\text{pos}$  and dimension  $i$ :

$$\text{PE}_{\{\text{pos}, i\}} = \sin\left(\frac{\text{pos}}{10000^{2i/d_{\text{model}}}}\right)$$

$$\text{PE}_{\{\text{pos}, i+1\}} = \cos\left(\frac{\text{pos}}{10000^{2i/d_{\text{model}}}}\right)$$

Create a positional encoding matrix of shape  $(\text{seq\_len}, d_{\text{model}})$ .

5. Add Token and Positional Embeddings: The input to the transformer is the element-wise sum of token embedding and positional Embedding

```
python ⌂ Copy ⌂ Edit
input_to_attention = token_embeddings + positional_encoding
```

Shape =  $(3, 512) \rightarrow \text{sequence length} \times \text{embedding dimension}$   
This is the final matrix passed into the transformer block.

#### PyTorch Style Code Skeleton:

```
python ⌂ Copy ⌂ Edit
import torch
```

### Decoder Block

Let's say we're decoding this: "Les trans"

1. Tokenization → Embeddings + Positional Encoding

```
Just like the encoder:
• Tokens → IDs → token embeddings  $x \in \mathbb{R}^{(\text{tgt\_seq\_len} \times d_{\text{model}})}$ 
• Add positional encoding (fixed or learned)

Let:
python ⌂ Copy ⌂ Edit
tgt_seq_len = 3
d_model = 512

You get:
text ⌂ Copy ⌂ Edit
x ∈ ℝ(3 × 512)
```

2. Masked Multi-Head Self Attention:

Allow each position to attend only to previous or same positions, not to future tokens.  
Why?

At training time, all tokens are known.

At inference, we generate tokens one by one — future tokens aren't available.

```
A triangular mask is applied:
text ⌂ Copy ⌂ Edit
Allowed:
Position 0 -> [0]
Position 1 -> [0, 1]
Position 2 -> [0, 1, 2]

So attention scores at position [1] are masked such that:
python ⌂ Copy ⌂ Edit
mask[[1][j]] = -inf if j > 1 else 0

This is applied before softmax:
python ⌂ Copy ⌂ Edit
attn = softmax(QK^T / sqrt_k + mask)
```

Let input  $x \in \mathbb{R}^{(3 \times 512)}$   
After Q, K, V projections (for all  $h$  heads):

- $Q, K, V \in \mathbb{R}^{(3 \times d_{\text{model}})}$
- For all heads → output shape =  $\mathbb{R}^{(3 \times 512)}$

3. Add & LayerNorm:

```
python ⌂ Copy ⌂ Edit
x1 = LayerNorm(x + MaskedSelfAttention(x))
```

4. Encoder-Decoder Cross-Attention

#### Goal:

Let the decoder attend to the encoder's output — the encoded source sequence.  
At each decoder position (e.g., French word being generated), we ask:  
"Which parts of the source sentence are relevant to me?"

```

import torch, mm as mm

# Dummy inputs
token_ids = torch.tensor([17, 56, 94]) # shape: (seq_len,)

# Embedding
embedding = nn.Embedding(10000, 512)
token_emb = embedding(token_ids) # (3, 512)

# Positional Encoding
def get_pos_enc(seq_len, d_model):
    pos = torch.arange(seq_len).unsqueeze(1)
    i = torch.arange(d_model).unsqueeze(0)
    angle_rates = 1 / torch.pow(10000, (i + (j // 2)) / d_model)
    angles = pos * angle_rates
    pe[:, :, :, 0] = torch.sin(angles)
    pe[:, :, :, 1] = torch.cos(angles)
    return pe

pos_enc = get_pos_enc(seq_len=3, d_model=512) # shape: (3, 512)

# Final Input to Attention
x = token_emb + pos_enc # (3, 512)

```

## 6. Multi-Head Self-Attention

At each position in the sequence, we want the token to understand the context it's in.

```

For example, in the sentence:
arkansas          # Copy 37 Edit
"Bank of the river"
"bank" should attend more to "river" not "money". So, we need a way for each token to "look at" all others
and decide who to pay attention to.

```

### A) Linear Projections for Q, K, V

We project  $x$  into queries ( $Q$ ), keys ( $K$ ), values ( $V$ ) using learnable matrices:

$$Q = xW^Q, \quad K = xW^K, \quad V = xW^V$$

- $W^Q, W^K, W^V \in \mathbb{R}^{(d\_model \times d\_k)}$  where  $d_k = d\_model / \text{num\_heads} = 64$

### B) Scaled Dot-Product Attention (per head)

For each head:

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

#### Intuition:

- $Q \cdot K^T$  gives a score: how much each token wants to pay attention to another.
- Sofmax turns this into a distribution over tokens.
- Multiplying by  $V$  tells the model: "What should I take from those tokens I'm attending to?"

Each token now becomes a weighted sum of all other tokens — but learned, context-aware, and dynamic for each head.

### C) Concatenate Heads & Final Linear Layer

```

MultiHead(x) = Concat(head_1, ..., head_h)W^O
Where  $W^O \in \mathbb{R}^{(d\_model \times d\_model)}$ 

```

#### Why multi-head?

Each attention head can focus on different relationships:

- One head may focus on subject-verb
- Another may track time relationships
- Another may resolve ambiguity ("bank" context)

Instead of one single attention operation, we have:

- Split  $d_{model}$  into  $h$  smaller dimensions (e.g., 512  $\rightarrow$  8 heads  $\times$  64)
- Run attention independently in each
- Concatenate outputs  $\rightarrow$  project back to  $d_{model}$

$$\text{MultiHead}(x) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

## 7. Add & LayerNorm

Apply residual connection and layer normalization:

$$x_1 = \text{LayerNorm}(x + \text{MultiHead}(x))$$

### Why residual connection?

The network can always choose to copy the input — which stabilizes training and helps gradient flow. It also allows the model to learn refinements, not entirely new representations.

### Why LayerNorm?

To normalize the input — this reduces internal covariate shift, stabilizes learning, and improves convergence. It ensures that all features have comparable scale and dynamics.

## 8. Feed-Forward Network (FFN)

Apply two linear layers with ReLU:

```

FFN(x) = \text{ReLU}(xW_1 + b_1)W_2 + b_2
Where:
•  $W_1 \in \mathbb{R}^{(d\_model \times d\_ff)}$ ,  $W_2 \in \mathbb{R}^{(d\_ff \times d\_model)}$ 

```

Typically,  $d_{ff}$  is  $4 \times d_{model}$

### Why FFN?

Attention mixes information across positions. FFN adds non-linearity and transformation power at each token independently.

Think of this as "refining" the representation for each token after it's looked at everyone else.

## 9. Add & Layer Norm (again)

Another residual connection:

$$x_2 = \text{LayerNorm}(x_1 + \text{FFN}(x_1))$$

## 10. Final Encoder output

### You stack N such blocks

$$\text{EncoderOutput} = \text{EncoderLayer}_N(\dots, \text{EncoderLayer}_2(\text{EncoderLayer}_1(x)))$$

Shape remains the same:  $\text{seq\_len} \times d_{model}$

But the content of each token's embedding now encodes contextualized information from all other tokens in the input.

```

# Q, K, V in Cross-Attention
# Q comes from the decoder input x1
# K, V come from encoder output enc_out
python
Q = x1 @ W^Q # shape: (tgt_seq_len x d_k)
K = enc_out @ W^K # shape: (src_seq_len x d_k)
V = enc_out @ W^V # shape: (src_seq_len x d_v)

# enc_out: E \mathbb{R}^{(src\_len \times d\_model)}
b. Dimensions:
• Q \in \mathbb{R}^{(tgt\_len \times d\_k)}
• K, V \in \mathbb{R}^{(src\_len \times d\_k)}
Then compute:
Attention(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V
This allows the decoder to attend over the source sentence at every position.

```

### Why not mask here?

Because:

- Source sentence is fully available, even during inference.
- So we do not apply masking in cross-attention.

## 5. Add & LayerNorm

```

python
x2 = LayerNorm(x1 + CrossAttention(x1, enc_out))

```

## 6. Feedforward Network

```

Same as in encoder:
python
FFN(x) = \text{ReLU}(xW_1 + b_1)W_2 + b_2
Where:
• W1 \in \mathbb{R}^{(512 \times 2048)}
• W2 \in \mathbb{R}^{(2048 \times 512)}

```

## 7. Add & LayerNorm

```

Final residual connection:
python
x3 = LayerNorm(x2 + FFN(x2))

```

## 8. Stack N Decoder Blocks

Just like the encoder, we repeat this  $N$  times.

Each decoder layer builds deeper representations, now jointly shaped by past target context + source context.

## 9. Final Output

```

Decoder output: E \mathbb{R}^{(tgt\_seq\_len \times d\_model)}
Pass this through a final linear layer and softmax:
python
logits = decoder_out @ W.Vocab.T # shape: (tgt\_len \times vocab\_size)
Used to compute probabilities of the next tokens.

```

## 2. KV Caching

KV caching is a critical optimization used during **inference (not training)** in **decoder-only** or **auto-regressive** transformer models like GPT. It helps speed up inference by avoiding redundant computation across time steps.

### Problem Before KV Caching (Without Cache)

When generating a sequence (e.g., a sentence), say "The cat sat", the transformer processes the entire prefix at every new token step.

1. To predict the first token, say "The" → compute  $\text{self\_attention}(\text{The})$
2. To predict "cat" → recompute attention over "The, cat"
3. To predict "sat" → recompute attention over "The, cat, sat"

This is redundant. Past key-value projections are recalculated again and again.

### KV Caching: Motivation

Avoid recomputing attention on tokens you've already processed. Instead, **cache the key and value tensors** from previous time steps.

#### How KV Caching Works

##### 1. Token Embedding

- Input: One token at a time (during inference), say token  $t_i$
- Output: Embedded representation  $x_i \in \mathbb{R}^{d_{\text{model}}}$

##### 2. Linear Projections to Q, K, V

For each input token  $x_i$ , compute:

- Query:  $Q_i = x_i \cdot W_Q$
- Key:  $K_i = x_i \cdot W_K$
- Value:  $V_i = x_i \cdot W_V$

Shapes:

- $Q_i$ : ( $\text{num\_heads}, 1, d_k$ )
- $K_i$ : ( $\text{num\_heads}, 1, d_k$ )
- $V_i$ : ( $\text{num\_heads}, 1, d_v$ )

##### 3. KV Caching Begins

- For the **first token**:
  - Compute  $K_1, V_1$ , and store them in the cache
  - $\text{cached\_K} = [K_1]$ ,  $\text{cached\_V} = [V_1]$
- For **subsequent tokens**:
  - Only compute  $K_i, V_i$  for the current token
  - **Append to cache**:  $\text{cached\_K} = [K_1, \dots, K_i]$ , same for  $\text{cached\_V}$

##### 4. Attention Computation

At each time step  $i$ :

- Use only the **new query**  $Q_i$
- Use all **cached keys/values**:  $\text{cached\_K} = [K_1, K_{i-1}], \text{cached\_V} = [V_1, V_{i-1}]$

```
Compute attention:
    Q_i · P_kit
Attention(Q_i, cached_K, cached_V) = softmax(Q_i / v_d_k) · cached_V
This gives the contextual representation for x_i.
```

##### 5. No Re-computation of Past K/V

Since the cached K/V are reused, the model avoids redoing linear projections for all past tokens. The expensive attention matrix grows only by 1 step per token.

##### 6. Caching Across All Layers

Each transformer block maintains its own KV cache. So if you have 24 layers:

- Each layer keeps  $\text{cached\_K}$ ,  $\text{cached\_V}$  per layer
- These are updated layer-wise per token during autoregressive decoding

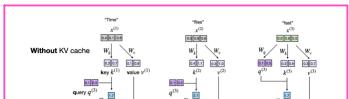
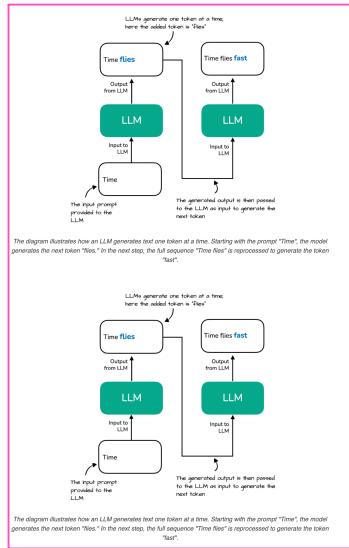
##### 7. Memory vs Speed Tradeoff

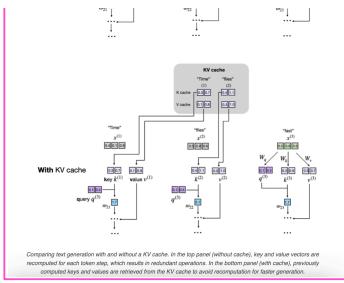
- Pros: Huge speedup during inference (especially for long sequences)
- Cons: Requires additional **memory** for storing cached K/V for every layer and head

##### 8. Why Not Cache Queries?

- Queries are computed only for the current token (during inference), and are **not reused**.
- Keys/Values are reused across all future time steps — only they are cached.

##### 9. Diagrammatic Representation:





Reference: <https://sebastianraschka.com/blog/2025/coding-the-kv-cache-in-lmns.html>

## 10. Size of KV Cache

$$L \times B \times N \times H \times S \times 2 \times 2$$

- L: Number of Transformer blocks
- B: Batch Size
- N: Number of Attention Heads
- H: Attention Head Size
- S: Context Length
- 2: No. of Caches per transformer block (K,V)
- 2: No. of Bytes per floating point Number (Assuming each parameter is 16 bit)

GPT-2: 36MB

GPT-3: 4.5GB

## 3. Multi Query Attention

**Multi-Query Attention (MQA)** is a variant of multi-head attention in transformers where:

- Each attention head has its own Query projection.
- All heads share the same Key and Value projections.

This design differs from standard multi-head attention (MHA), where each head has its own Q, K, and V projections.

### 1. Input Embeddings

- o Suppose you have an input sequence of shape (batch\_size, seq\_len, d\_model).
- o Let's assume d\_model = 768 and num\_heads = 12.

### 2. Linear Projections

- o In standard MHA, you create separate projection matrices for Q, K, and V per head.
  - So you get  $Q = W_Q \cdot X$ ,  $K = W_K \cdot X$ ,  $V = W_V \cdot X$ , where all are shaped (batch\_size, seq\_len, d\_k) for each head.
  - Resulting in 12 Q, 12 K, and 12 V for 12 heads.

### 3. KV Caching in Decoding

- o In autoregressive decoding (e.g., GPT), to generate one token at a time efficiently, we cache past K and V values for reuse.
- o In standard MHA, each head has its own K/V → you must store (num\_heads, seq\_len, d\_k) for both K and V.
- o In MQA, all heads share the same K and V → you only store one set of (seq\_len, d\_k) for K and V.

### 4. Memory Savings

- o Significant reduction in KV cache size during inference.
  - From 12x to 1x for K and V.
  - Especially important for long-context generation and edge devices.

### 5. Faster Generation

- o Less memory → better cache locality.
- o Smaller key-value caches → faster memory reads.
- o Improved latency in autoregressive decoding (e.g., LLaMA-2 Chat uses MQA for this reason).

### 6. Loss of Head Diversity

- o In standard MHA, each head can learn different semantic attention patterns via separate KV projections.
- o MQA limits this by forcing all heads to use the same K and V → reduces model expressiveness.
- o Heads may collapse into similar attention maps, since only Q differs.

### 7. Reduced Representational Power

- o With shared KV, model has less flexibility to attend to diverse contextual signals.
- o May hurt fine-grained attention, especially in tasks needing nuanced alignment (e.g., translation, QA).

### 8. Accuracy Drop in Some Tasks

- o Empirically, MQA can perform worse than MHA on complex reasoning or multilingual tasks.
- o Works well for text generation, but not always for other modalities (vision, multimodal) or for short context tasks.

### 9. Trade-off Between Speed and Accuracy

- o MQA trades speed and memory efficiency for a small drop in quality.
- o Acceptable in large models like GPT-3.5, LLaMA-2~7B+, but may not suit small models or low-resource training.

## 10. Size of KV Cache

$$L \times B \times H \times S \times 2 \times 2$$

L: Number of Transformer blocks

B: Batch Size

H: Attention Head Size

S: Context Length

2: No. of Caches per transformer block (K,V)

2: No. of Bytes per floating point Number (Assuming each parameter is 16 bit)

GPT-3 with MHA : 4.5GB  
GPT-3 with MQA: 48MB

## 4. Grouped Query Attention

Grouped Query Attention is a hybrid between Multi-Head Attention (MHA) and Multi-Query Attention (MQA) where:

- Each group of attention heads shares a single Key and Value projection.
- But each head still gets its own Query projection, like in MHA.

This provides a balance between MQA's memory efficiency and MHA's representational power.

## 1. Input Embeddings

- Suppose you have an input sequence with shape (batch\_size, seq\_len, d\_model).
- Let's assume  $d_{model} = 768$  and  $num\_heads = 12$ .
- Let's say you group the 12 heads into 4 groups, so each group has 3 heads.

## 2. Linear Projections

- In standard MHA:
  - Each head has its own Q, K, and V projection.
  - You get 12 Qs, 12 Ks, and 12 Vs.
- In GQA (for 4 groups and 12 heads):
  - You have:
    - 12 different  $W_Q$  matrices → one Query projection per head.
    - Only 4 different  $W_K$  and 4 different  $W_V$  matrices → one shared Key and Value per group.
  - So:
    - $Q_1, \dots, Q_{12}$  (1 per head)
    - $K_1, \dots, K_4$  (1 per group)
    - $V_1, \dots, V_4$  (1 per group)
  - Each group of 3 heads shares the same K and V but uses their own Q.

## 3. KV Caching in Decoding

- During autoregressive decoding:
  - In MHA: You cache 12 separate K/V → expensive.
  - In MQA: You cache only 1 K/V → very efficient.
  - In GQA (with 4 groups): You cache 4 sets of K/V → a middle ground.
- KV cache shape becomes ( $num\_groups, seq\_len, d_k$ ) instead of ( $num\_heads, seq\_len, d_k$ ).

## 4. Memory Savings

- Moderate reduction in KV cache size compared to full MHA.
  - If  $num\_heads = 12$  and  $num\_groups = 4$ , cache size drops by 3x.
- This is useful for long-sequence generation and deployment on limited hardware.

## 5. Faster Generation

- Fewer KV tensors → better memory locality.
- Smaller cache → faster memory access during token generation.
- Improves inference speed, though less so than MQA.

## 6. Preserves Partial Head Diversity

- Unlike MQA, not all heads share the same K/V.
- Each group can attend to different contexts via their own K/V.
- Retains some of the expressiveness of standard MHA.
- Heads within a group may correlate more, but inter-group diversity remains.

## 7. Better Representational Power than MQA

- Because K/V are shared per group, not globally:
  - More fine-grained attention than MQA.
  - Still less expressive than full MHA, but better suited for complex reasoning tasks than MQA.

## 8. Improved Accuracy over MQA

- Empirical studies show GQA performs better than MQA on tasks like:
  - Multilingual understanding
  - Question answering
  - Token classification
- Widely used in models like LLaMA-3, Mistral, and Gemma.

## 9. Trade-off: Accuracy vs. Efficiency (Sweet Spot)

- GQA offers a balance:
  - Much faster and more memory-efficient than MHA.
  - Retains much more accuracy than MQA.
- Ideal for models where:
  - Scalability is key (e.g., LLMs with long contexts)
  - Edge deployments or limited compute environments matter
  - High-quality outputs are still necessary.

## 10. Size of KV Cache

$$L \times B \times G \times H \times S \times 2 \times 2$$

L: Number of Transformer blocks

B: Batch Size

G: No. of Groups

H: Attention Head Size

S: Context Length

2: No. of Caches per transformer block (K,V)

2: No. of Bytes per floating point Number (Assuming each parameter is 16 bit)

Bottom line:

GQA ≈ "Goldilocks zone" — not as expensive as MHA, not as lossy as MQA. A solid compromise used in modern LLMs.

## 5. Multi Head Latent Attention

Diagram from paper - DeepSeek-V2:

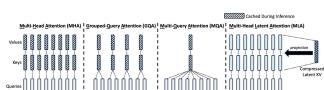


Figure 3 | Simplified illustration of Multi-Head Attention (MHA), Grouped-Query Attention (GQA), Multi-Query Attention (MQA), and Multi-Head Latent Attention (MLA). Through jointly compressing the keys and values into a latent vector, MLA significantly reduces the KV cache during inference.

Attention Mechanism Type	Number of Unique Keys/Values	Resources	Performance / Context Understanding
Multi-Head Attention (MHA)	Each head has its own key and value	Largest	Best
Multi-Query Attention (MQA)	All heads share same key and value	Smallest	Worst
Grouped Query Attention (GQA)	Heads are divided in G groups ( each group share key and value)	Medium	Medium

Can we get best of both worlds? KV cache size and performance?

Size of KV cache  
 $L \times B \times G \times H \times S \times 2 \times 2$

Good Language Performance Model  
 For this, All heads should have different values in K, V matrix.

This needs to be resolved.

How do we achieve this?

What if we don't have to cache the Key and Values Separately?  
What if we cache only one matrix?  
What if this matrix has lesser dimensions than NxD?

To achieve this let's start by projecting this matrix into a latent space.

$X = (4 \times 8)$



$W_{dKv} = (8 \times 4)$

$C_{kv} = 4 \times 4$

This is our latent matrix which we will cache.

What do we do next?

$X(4 \times 8)$



$K()$

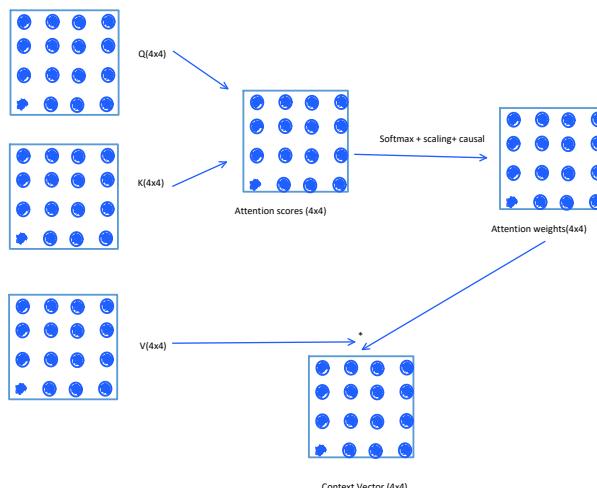
$W_{dkv} = (8 \times 4)$

= Att

$C_{kv} (4 \times 4) =$

$* W_{uk}^T (4 \times 4)$

$* W_{uv}^T (4 \times 4)$



Context Vector (4x4)

How does adding this latent matrix help?

1.  $Q = X * W_Q$

2.  $C_{kv} = X * W_{dkv}$

3.  $K = C_{kv} * W_{uk} = X * W_{dkv} * W_{uk}$

4.  $V = C_{kv} * W_{uv} = X * W_{dkv} * W_{uv}$

The absorption trick

$\text{Attention scores} = Q * K^T$

$= (X * W_Q) * (W_{uk}^T * W_{dkv}^T * X^T)$

$= (X * W_Q * W_{uk}^T) * (X * W_{dkv})^T$

$= (X * W_Q * W_{uk}^T) * (C_{kv})^T$

This is fixed during training.

This needs to be cached

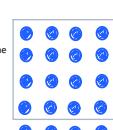
Context vector = Attention weights \* V

$= (Q * K^T) * (X * W_{dkv} * W_{uv}) * W_V$

$= \text{Attention scores} * C_{kv} * (W_{uv} * W_V)$

Already Computed | Cached | fixed at training

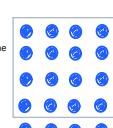
What happens when the new token comes in?



1. Create Query (Q)  
 $Q = X_{1x8} * W_Q_{8x4} * W_{uk}_{4x4}$

2. Compute KV vector

$= X_{1x8} * W_{dkv}_{8x4} \longrightarrow \text{Append this to KV cache}$



3. Now we multiply updated KV cache with absorbed query vector

Attention scores =  $Q_{1x4} * \text{Updated KV}_{5x4}^T$

Attention score will have dimension (1x5)

4. Compute attention weights for the same

Attention weights will also have dimension 1x5

5. Calculate Value matrix

Value =  $\text{Updated KV}_{5x4}^T * W_{uv}_{4x4}$

Value matrix will have dimension (5x4)

#### 6. Calculate context vector

Context Vector<sub>1..n</sub> = Attention Weights<sub>1..n</sub> \* V<sub>E..n</sub>

Did we actually solved the two problems we started with?

**Size of KV cache**  
 $L \times B \times N \times H \times S \times 2 \times 2$

**Good Language Performance Model**  
For this, All heads should have different values in K, V matrix

This Needs to be reduced.

Size in KV cache after MLA

$$L \times B \times dL \times S \times 2 \quad \longrightarrow \quad \text{Reduction in size: } \frac{2 \times N \times H}{dL}$$

For Deenseek:

$$= \frac{2 * 128 * 128}{576} \approx 57 \text{ times}$$

= 400GB to 6 GB

## 6. Positional Embedding

- A) Integer Positional Embedding
  - B) Sinusoidal Positional Embedding
  - C) Rotary Positional Embedding

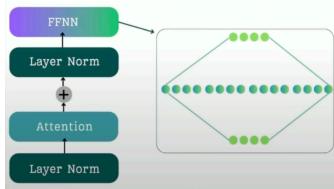
## 7. Multi Head Latent Attention With Rotary Positional Embedding

## 2. Mixture of Experts

Another Innovative thing deepseek did is creative usage of Mixture of Experts Technique. In traditional neural Networks standard feedforward neural network is used.

## 1. Traditional Mixture of Experts

It's not a new idea, Mixture of expert was introduced in 1991 in neural computation.



As shown above, Feedforward network has one hidden layer with 4 times number of units as the input layer.

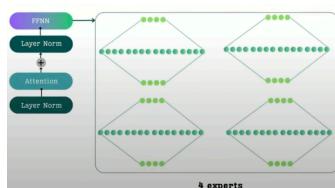
The output has same number of units as in input layer

### Parameters:

Answers 1-360 • The McGraw-Hill Companies

Same number of parameters

As parameters are huge this increases training and inference time. Mixture of experts come to rescue here.

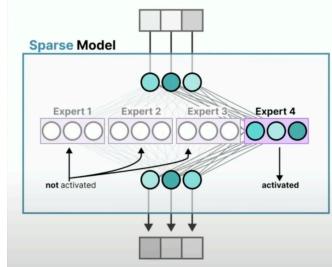


In figure we have shown 4 experts, these are per transformer block.  
In mixture of experts model we have multiple of such networks called experts in transformer block. So how does it help than FFN (dense).

In mixture of experts model we have multiple of such networks called experts in parallel.

In a dense model, input token passes through all the parameters, (all layers and neurons)

In mixture of input we have multiple experts but only small subset of them are activated for any given input



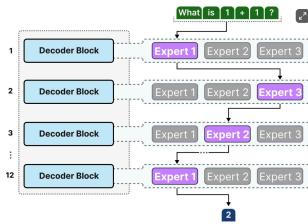
What do these experts learn? <https://arxiv.org/pdf/2202.08908.pdf>

	sen sau road miss place smuggling tuning slogan seo falling designed based smuggling submitted develop her over her know dark upper dark outer center very bright bright bright over open your dark black
Visual descriptions color, spatial position	Layer 0
Proper name	Layer 1 A Mart Or Mart Kone Mod Mod Cor Tri Ca Ma R Mart Lorraine Colin Ken Sam Ken Ce Angel A Dong Nguon Giai Phuoc Giai Giai Giai

Counting and numbers  
written and numerical form

Layer 1  
after 37 09 00 27 11 Seven 25 4, 54 Two dead we  
Some 2012 who we for lower each

Now the token travel to step?

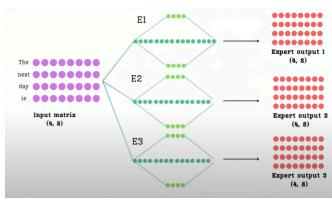


Now, We'll go through hands on example to see how MOE actually works:

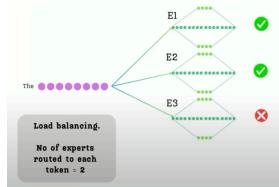
Step 1 We start with input matrix of size 4x8 (4 token with embedding dimension 8)



We are going to pass this input matrix to MOE with 3 experts. Note that each expert is nothing but simple FNN.



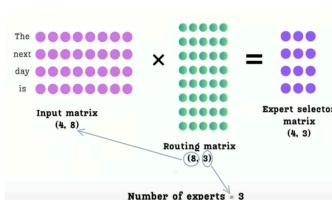
Step 2 When we pass input matrix through these 3 experts, we get 3 expert outputs matrices each of 4x8.  
The challenge is to figure out how to merge these to get single resulting matrix of 4x8



Step 3 To do the merging, we use a method called load balancing.  
Every token will be sent to only selected number of experts. This is called load balancing.  
As shown in figure, lets say that we decide to route token only via 2 experts.

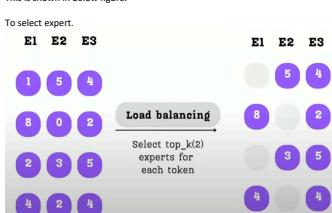
Step 4 Once we have decided how many experts will be assigned to each token, the next question is  
"How much weightage to give to each expert?"  
"Which Expert to select?"

This is decided using Routing mechanism.

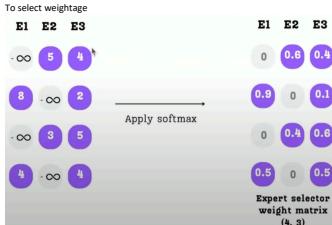


As shown in figure above, we multiply the input matrix 4x8 with a routing matrix 8x3 and this leads to a expert selector matrix 4x3  
Note that every column of expert selector corresponds to each expert.  
Since we only decide to keep 2 experts, for each token, we take expert selector matrix and make sure each row has only 2 active experts.

This is shown in below figure.

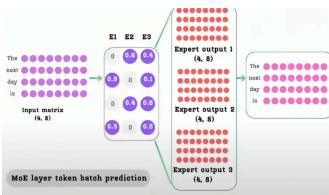


To select weightage



Step 5

Once we have expert selector weight matrix. We can use then to merge expert matrices.



There are two major issues with Mixture of Experts. We'll discuss those in next two steps.

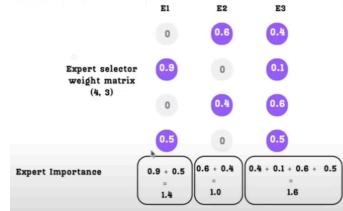
#### Step 6: Auxiliary loss version 1

In MoE models, the routing mechanism selects a subset of experts for each input. If some expert is chosen too often while others remain unutilized, this creates inefficient learning and performance bottlenecks.

An auxiliary loss term is added to main training loss to penalize imbalanced expert selection, pushing the routing function towards more uniform distribution. To calculate auxiliary loss, we first start with expert selector weight matrix which consists of experts assigned to every token, and probabilities assigned to every expert.

Each column of expert selector weight matrix gives us probability assigned to each expert. Based on this we calculate term called expert importance.

Expert importance is just sum of all probabilities for that expert.



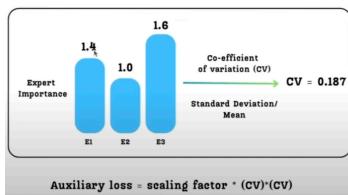
We want a balanced MoE. What does balance mean?

We want roughly same number of tokens routed to each expert. How to quantify it? We just saw above. Expert importance.

We want to penalize model if expert importance has more variance. This would mean some experts are more important than others.

The coefficient variation is given by formula:

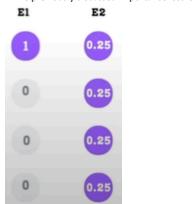
Coefficient Variation = Standard Deviation / Mean



We add this auxiliary loss to l1m loss.

#### Step 7: Load Balancing

While previously discussed importance loss is useful. Assigning equal importance to experts does not necessarily lead to uniform token routing.



Here both experts adds up to 1 and hence have same importance but expert 2 is getting selected way more than expert 1.

We need another way: load balancing.

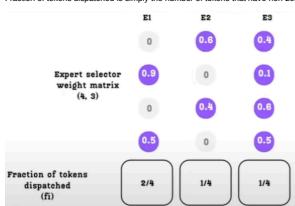
The first step of load balancing is to calculate probability that the router choose given expert.

We can calculate it by dividing expert importance by number of tokens

- Probability that the router will choose Expert 1 ( $p_1$ ) =  $1.4/4 \approx 0.35$
- Probability that the router will choose Expert 2 ( $p_2$ ) =  $1/4 \approx 0.25$
- Probability that the router will choose Expert 3 ( $p_3$ ) =  $1.6/4 \approx 0.4$

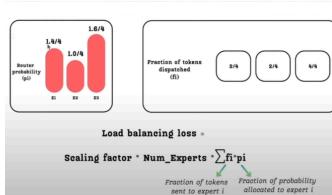
The next step is to calculate the fraction of tokens dispatched to each expert.

Again, this can be calculated from the expert selector weight matrix. Fraction of tokens dispatched is simply the number of tokens that have non-zero activation for each expert divided by total tokens (in this case = 4).



- The fraction of tokens dispatched to expert 1 ( $f_1$ ) =  $2/4 = 0.5$
- The fraction of tokens dispatched to expert 2 ( $f_2$ ) =  $1/4 = 0.25$
- The fraction of tokens dispatched to expert 3 ( $f_3$ ) =  $1/4 = 0.25$

$$\text{Load balancing loss} = \text{Scaling factor} * \text{Number of experts} * \sum_{i=1}^{\text{num-experts}} f_i p_i$$



$p_i$  is the fraction of total router probability (importance) given to expert i.

$f_i$  is the fraction of tokens dispatched to expert i.

By minimizing the auxiliary loss, you mathematically enforce the model to distribute tokens proportionally to how much each expert is valued or "trusted."

resulting in more efficient and balanced use of the MoE architecture.

#### Step 8: Capacity Factor

There is one more issue which a MoE model faces: expert imbalance.

Expert overload essentially means some experts receiving disproportionately high number of tokens, while others not receiving any.

Although we have seen another ways to handle this, this is also a guardrail we can enforce.

Expert capacity is the maximum number of tokens that a single expert can process in a given batch.

Expert capacity = Tokens per batch / Number of experts \* Capacity Factor

Tokens per batch = batch size \* sequence length \* top\_k

[DeepSeek: Enabling Auxiliary Loss-Free Load Balancing](#)

Load balancing loss = Scaling factor \* Number of experts \*  $\sum_{i=1}^n f_i p_i$

This is the load balancing loss we saw earlier.

The use of this loss helps maintaining load balance in expert models but it also acts as a regularization term that interferes with language modeling.

If the scaling factor is too small or absent, it leads to poor balance and routing collapse, reducing model efficiency and limiting expert utilization.

A large scaling factor ensures load balance but degrades overall model performance.

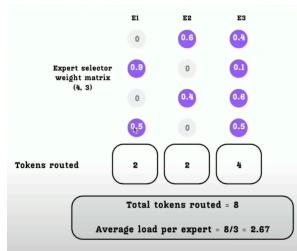
[Is load balancing a bottleneck problem?](#) balancing experts effectively without compromising training quality.

Can we do something better? Answer is always, Yes.

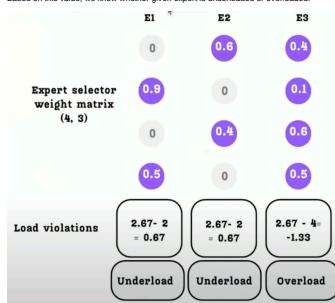
To resolve this problem, Deepseek implemented Loss-Free Balancing, a method that directly enforces load balance without introducing additional gradients beyond those from the language modeling loss.

Here is how its done:

Based on the total tokens routed, we calculate the average load per token.



Based on this value, we know whether given expert is underloaded or overloaded.



From above figure expert 1 and 2 are underloaded and expert 3 is overloaded.

Now, we look at load violation error = average load per expert - actual load per expert

In the next step, we introduce bias.

For each expert, we introduce a bias term b.

The bias is initialized as 0 for each expert.

Here is how the bias terms are updated for each expert:

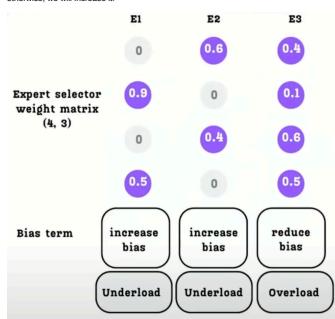
$$b_i = b_i + u * \text{sign}(\text{load violation error})$$

Here, u is predefined constant.

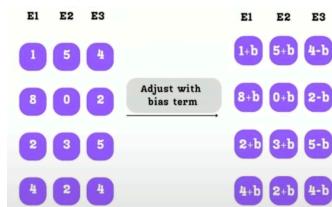
Note that if an expert is underloaded, its load violation error is positive.

If an expert is overloaded, its load violation error is negative.

The above formula ensures that if an expert has a heavy load on the previous batch, we will reduce its bias. otherwise, we will increase it.



Here is how its done



Since expert 1 and 2 are underloaded, adding these bias terms will make sure values of the expert 1 and expert 2 increase.

Through the dynamic adjustment for the biases, Deepseek achieved good expert load balance, without directly introducing noisy gradients into the model like what happens with the load balancing loss.

Deepseek showed that auxiliary loss free load balancing achieves both better performance and better load balance compared with traditional load balancing.

[Deepseek: Enabling Auxiliary Loss-Free Load Balancing](#)

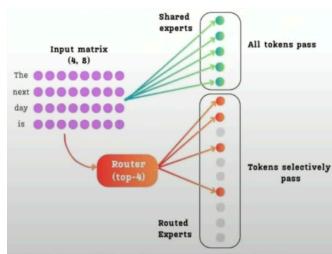
This architecture manifests two potential issues: (1) **Knowledge Hybirdity**: existing MoE practices often employ a limited number of experts (e.g., 8 or 16), and thus tokens assigned to a specific expert will be likely to cover diverse knowledge. Consequently, the designated expert will intend to assemble vastly different types of knowledge in its parameters, which are hard to utilize simultaneously. (2) **Knowledge Redundancy**: tokens assigned to different experts may require common knowledge. As a result, multiple experts may converge in acquiring shared knowledge in their respective parameters, thereby leading to redundancy in expert parameters. These issues collectively hinder the expert specialization in existing MoE practices, preventing them from reaching the theoretical upper-bound performance of MoE models.

Deepseek modified the existing MoE architecture by introducing the idea of shared experts.

This approach divides experts into two groups:

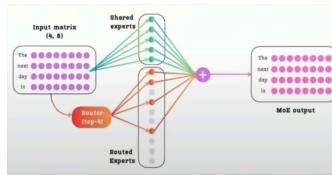
- Shared Experts: Experts that process every token, regardless of routing.
- Routed Experts: Experts that handle tokens selectively, based on the usual routing strategy.

While routed experts continue to be assigned tokens according to standard MoE procedures, shared experts are always active and handle every token in the batch.



Typically, the number of shared experts should be smaller than the number of routed experts.

If there are too many shared experts, the model becomes less sparse, reducing the efficiency benefits of using an MoE architecture.



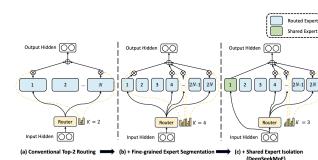
Deepseek Innovation 3: Fine-grained expert segmentation

If the number of experts is small, each expert is forced to learn a wide variety of knowledge types, reducing its specialization.

This leads to inefficient learning because the expert parameters need to generalize across too many diverse inputs.

In fine-grained expert segmentation, each large expert FFN (Feed-Forward Network) is split into smaller experts by reducing the hidden dimension of the FFN by a factor of 1/4.

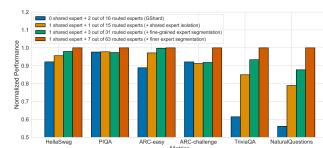
To compensate for the reduced capacity per expert, the model activates 4 times more experts per token while keeping the same total computational cost.



It is noteworthy that across these three architectures, the number of expert parameters and computational costs remain constant.

Some results from Deepseek-MoE paper:

Metric	# Shot	Dens	Hash Layer	Switch	GShard	DeepSeekMoE
# Total Params	N/A	0.2B	2.0B	2.0B	2.0B	DeepSeekMoE
# Activated Params	N/A	0.2B	0.2B	0.3B	0.3B	
tokens per 2K Tokens	N/A	1.0B	1.0B	1.0B	1.0B	
# Trained Tokens	N/A	100B	100B	100B	100B	
Pile (Loss)	N/A	2.060	1.932	1.881	1.867	<b>1.808</b>
HellSwoop (Acc.)	0-shot	38.8	46.2	49.1	50.5	<b>54.8</b>
PiQA (Acc.)	0-shot	66.8	68.4	70.5	70.7	<b>72.3</b>
ARC-easy (Acc.)	0-shot	41.0	45.3	45.9	43.9	<b>49.4</b>
ARC-challenge (Acc.)	0-shot	25.0	25.2	31.9	31.9	<b>34.3</b>
RACE-middle (Acc.)	5-shot	38.8	38.8	43.6	42.1	<b>44.0</b>
RACE-high (Acc.)	5-shot	29.0	30.0	30.9	30.4	<b>31.7</b>
HumanEval (Pass@R)	0-shot	0.0	1.2	2.4	3.7	<b>4.9</b>
MBPP (Pass@1)	3-shot	0.2	0.6	0.4	0.2	<b>2.2</b>
TriviaQA (EM)	5-shot	4.9	6.5	8.9	10.2	<b>16.6</b>
NaturalQuestions (EM)	5-shot	1.4	1.4	2.5	3.2	<b>5.7</b>



### 3. Multi Token Prediction

One of the major innovations in the Deepseek architecture is multi-token prediction.

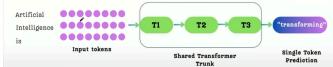
Multi-Token Prediction was a technique already introduced by Meta in their 2024 paper titled "Better & Faster Large Language Models via Multi-token Prediction".

(<https://arxiv.org/pdf/2404.19737.pdf>)

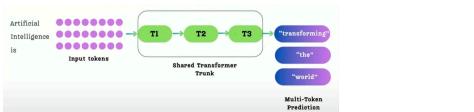
Deepseek implemented this same technique with slight modifications.

Let us understand what multi-token prediction is and how it was implemented by DeepSeek.

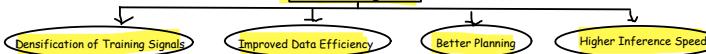
#### Single Token Prediction



#### Multi Token Prediction



#### Advantages



##### 1. Densification of training signals

MTP provides richer (long range understanding) and denser training signals than single-token prediction.

Traditional next-token training only guides the model to predict a single immediate token.

MTP, however, instructs the model to simultaneously predict multiple future tokens, generating more informative gradient signals per training example.

##### 2. Improved Data Efficiency

MTP-trained models achieved better results on standard benchmarks like HumanEval and MBPP with the same amount of training data, solving about 15% more code problems on average.

Training data	Vocabulary	n	MBPP		HumanEval		APPS/Intro	
			0@1	0@10	0@100	0@10	0@100	0@1
1	19.3	42.4	64.7	18.1	28.2	47.8	0.1	0.5

31.3B bytes (0.5 epochs)	bytes	8	<b>32.3</b>	<b>50.0</b>	<b>69.6</b>	<b>21.8</b>	<b>34.1</b>	<b>57.9</b>	<b>1.2</b>	<b>5.7</b>	<b>14.0</b>
		16	28.6	47.1	68.0	20.4	32.7	54.3	1.0	5.0	12.9
		32	30.0	53.8	76.7	22.8	36.4	62.0	2.8	8.8	21.8
2000 tokens (0.5 epochs)	32k tokens	1	30.0	53.8	76.7	22.8	36.4	62.0	2.8	8.8	17.4
		2	30.3	55.1	76.2	22.2	38.5	62.6	2.1	9.0	21.7
		4	23.8	<b>55.9</b>	<b>76.7</b>	<b>24.0</b>	<b>40.1</b>	<b>66.1</b>	1.6	7.1	19.9
		8	23.8	55.9	76.7	24.0	40.1	66.1	1.6	7.1	19.9
		16	20.7	52.2	73.4	20.0	36.6	59.6	3.5	10.4	22.1
17 tokens (4 epochs)	32k tokens	1	40.7	65.4	83.4	31.7	57.6	83.0	5.4	<b>17.8</b>	34.1
		4	<b>43.1</b>	65.9	83.7	31.6	57.3	<b>86.2</b>	4.3	15.6	33.7

### 3. Better Planning

Greater importance to choice points.

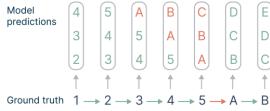


Figure 9: **Multi-token prediction loss assigns higher implicit weights to consequential tokens.** Shown is a sequence in which all transitions except “5 → A” are easy to predict, alongside the corresponding prediction targets in 3-token prediction. Since the consequences of the difficult transition “5 → A” are likewise hard to predict, this transition receives a higher implicit weight in the overall loss via its correlates “3 → A”, ..., “5 → C”.

### 4. Higher inference speed

Up to 3x faster inference speed.

### 3.2. Faster inference

We implement greedy self-speculative decoding (Stern et al., 2018) with heterogeneous batch sizes using xFormers (Lefauveux et al., 2022) and measure decoding speeds of our best 4-token prediction model with 7B parameters on completing prompts taken from a test dataset of code and natural language (Table S2) not seen during training. We observe a speedup of 3.0× on code with an average of 2.5 accepted tokens out of 3 suggestions on code, and of

Deepseek only used it during pretraining and not in inference.

### How Deepseek Rewrite Multi Token Prediction?



Inspired by Gaoček et al. (2024), we investigate and set a Multi-Token Prediction (MTP) objective for Deepseek-V3, which extends the prediction scope to multiple future tokens at each position. On the one hand, an MTP objective densifies the training signals and may improve data efficiency. On the other hand, MTP may enable the model to pre-plan its representations for the next tokens. In this section, we first introduce the MTP objective and then discuss the MTP losses from Gaoček et al. (2024), which parallelly predicts  $\theta$  additional tokens using independent output heads; we sequentially predict additional tokens and keep the complete causal chain at each prediction depth. We introduce the details of our MTP implementation in this section.

