# DeepSeek-V3:
# Multihead Latent Attention (MLA) & DeepSeekMoE

Minwu Kim, Safal Shrestha

deepseek

**DeepSeek-V3 Technical Report**

DeepSeek-AI

research@deepseek.com

## Abstract

We present DeepSeek-V3, a strong Mixture-of-Experts (MoE) language model with 671B total parameters with 37B activated for each token. To achieve efficient inference and cost-effective training, DeepSeek-V3 adopts Multi-head Latent Attention (MLA) and DeepSeekMoE architectures, which were thoroughly validated in DeepSeek-V2. Furthermore, DeepSeek-V3 pioneers an auxiliary-loss-free strategy for load balancing and sets a multi-token prediction training objective for stronger performance. We pre-train DeepSeek-V3 on 14.8 trillion diverse and high-quality tokens, followed by Supervised Fine-Tuning and Reinforcement Learning stages to fully harness its capabilities. Comprehensive evaluations reveal that DeepSeek-V3 outperforms other open-source models and achieves performance comparable to leading closed-source models. Despite its excellent performance, DeepSeek-V3 requires only 2.788M H800 GPU hours for its full training. In addition, its training process is remarkably stable. Throughout the entire training process, we did not experience any irrecoverable loss spikes or perform any rollbacks. The model checkpoints are available at https://github.com/deepseek-ai/DeepSeek-V3.
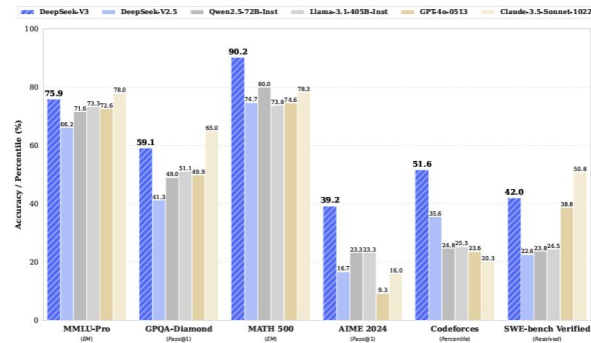
Figure 1 | Benchmark performance of DeepSeek-V3 and its counterparts.

# DeepSeek-V3

- 671B-parameter model with 37B active parameters during inference
- Performance on par with GPT-4o
- Backbone of DeepSeek-R1
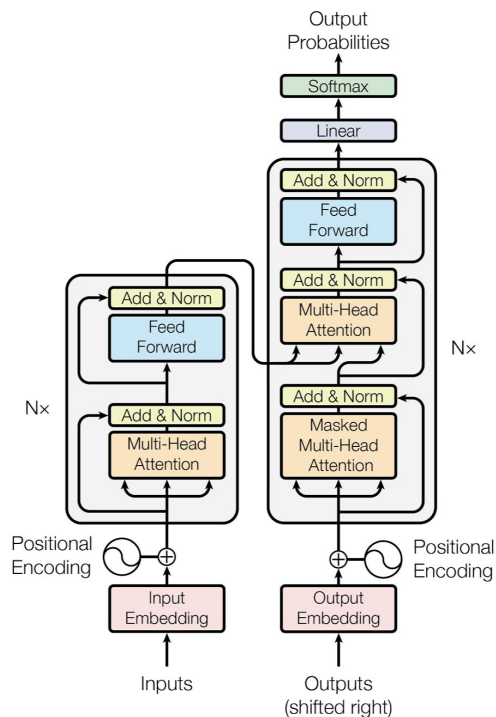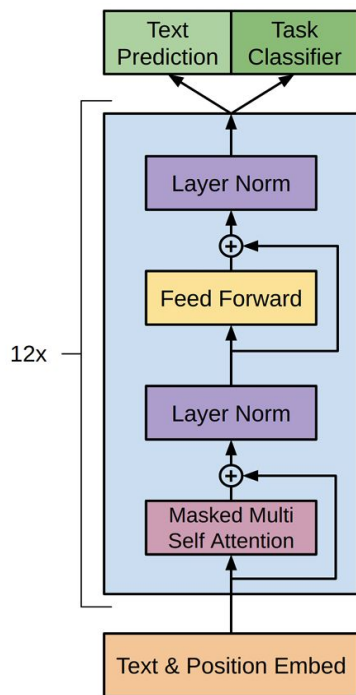- Optimized for efficient inference and cost-effective training

# DeepSeek-V3

Efficiency achieved by:
- Multi-head Latent Attention
- DeepSeekMoE

# Overview

1.  Recap for transformer
2.  Multi–head Latent Attention (MLA)
    a.  Self–attention
    b.  KV cache
    c.  MLA
3.  DeepSeekMoE
    a.  Feedforward component
    b.  Mixture of Experts (MoE)
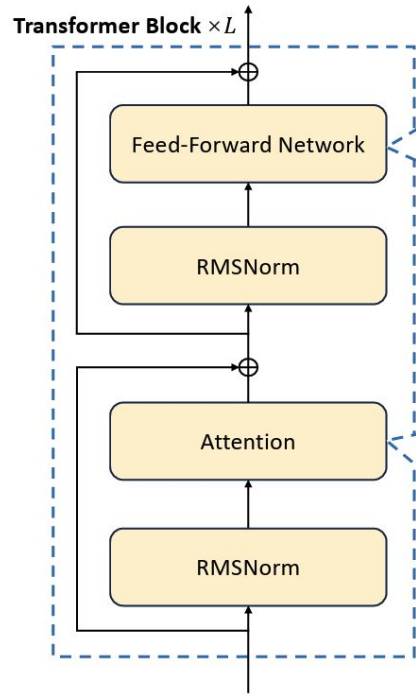    c.  DeepSeekMoE

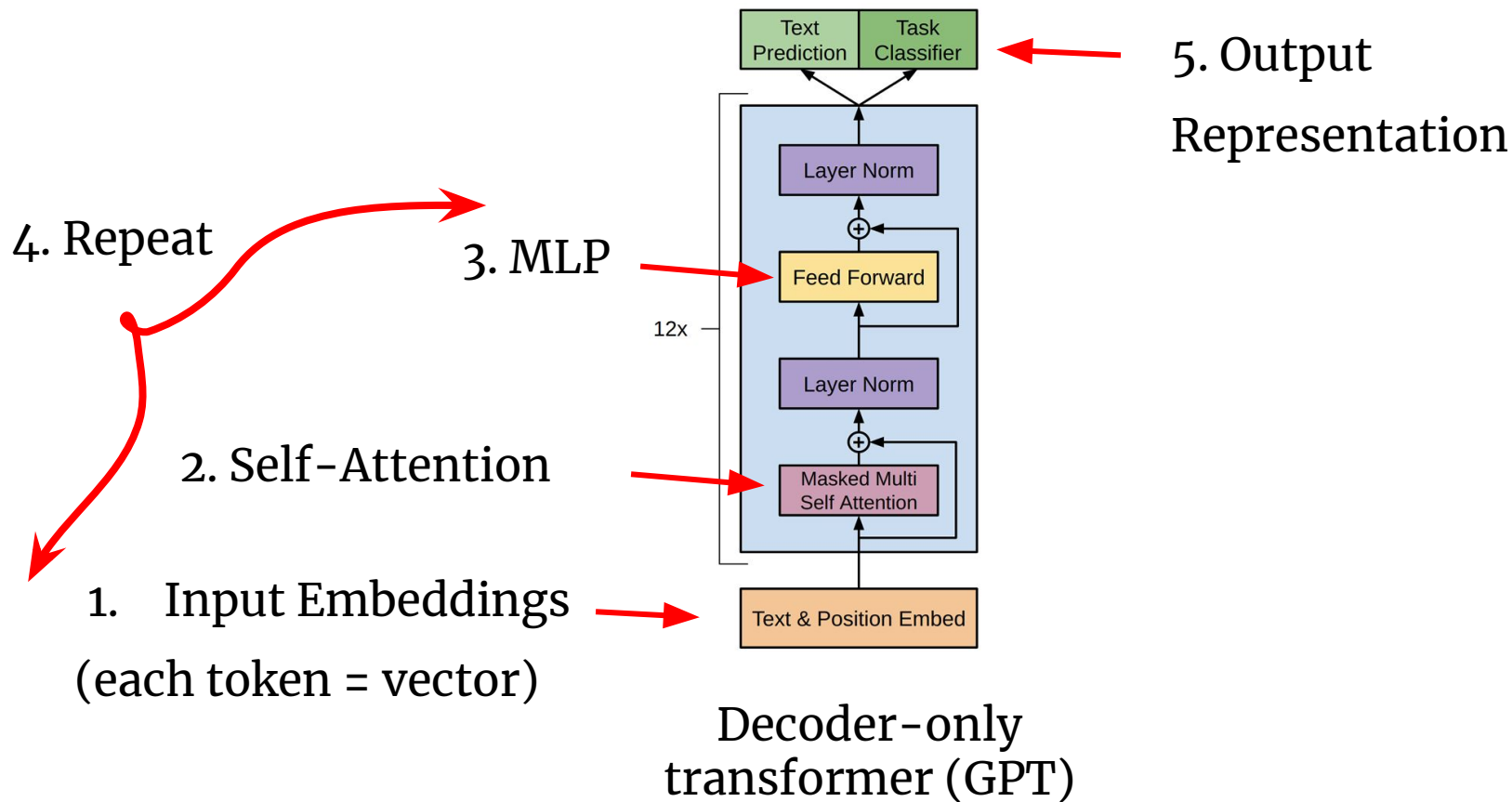# Preliminaries – Transformer Recap



Encoder-decoder transformer

Decoder-only transformer (GPT)

Decoder-only transformer (V3)
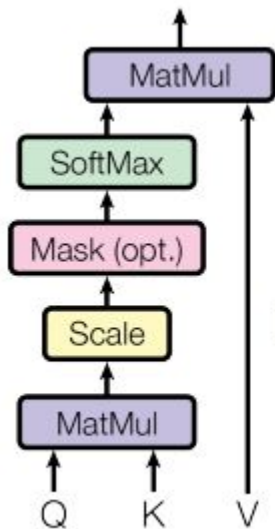
# Preliminaries – Transformer Recap



5. Output Representation

4. Repeat

3. MLP

2. Self-Attention

1. Input Embeddings
(each token = vector)

Text Prediction

Task Classifier

Layer Norm

Feed Forward

12x

Layer Norm

Masked Multi Self Attention

Text & Position Embed

Decoder-only transformer (GPT)

# Multi-head Latent Attention (MLA)

## Preliminary 1: Self-attention



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

# Multi-head Latent Attention (MLA)

## Preliminary 1: Self-attention

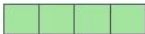$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

# Multi-head Latent Attention (MLA)

## Preliminary 1: Self-attention

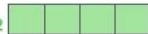$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

# Multi-head Latent Attention (MLA)

## Preliminary 2: Multi-head Attention

# Multi-head Latent Attention (MLA)

## Preliminary 3: Key-Value (KV) Caching



**Step 1**

Without cache

Q — Query Token 1 — (1, emb_size)

x

$K^T$ — Key Token 1 — (emb_size, 1)

=

$QK^T$ — $Q_1.K_1$ — (1, 1)

V — Value Token 1 — (1, emb_size)

x

=

Attention — Token 1 — (1, emb_size)

With cache

Q — Query Token 1 — (1, emb_size)

x

$K^T$ — Key Token 1 — (emb_size, 1)

=

$QK^T$ — $Q_1.K_1$ — (1, 1)

V — Value Token 1 — (1, emb_size)

x

=

Attention — Token 1 — (1, emb_size)

Values that will be masked   Values that will be taken from cache

# Multi-head Latent Attention (MLA)

## Preliminary 3: Key-Value (KV) Caching

# Multi-head Latent Attention (MLA)

## Preliminary 3: Key-Value (KV) Caching

# Multi-head Latent Attention (MLA)

## Preliminary 3: Key-Value (KV) Caching

# Multi-head Latent Attention (MLA)

## Preliminary 3: Key-Value (KV) Caching

# Multi-head Latent Attention (MLA)

**Problem: High Memory Consumption**

- As the context window grows, the KV cache scales its memory usage significantly, potentially straining system resources.

Motivation:

Is it possible to optimize or reduce the KV cache during inference without compromising performance?

# Multi-head Latent Attention (MLA)

Can we reduce the KV cache during inference?

- 8 attention heads -> Cache information of 8 KV pairs

    Idea: Jointly compress the keys and values into a latent vector and project them back to original vector space.

# Multi-head Latent Attention (MLA)



$W^{DKV}$

$\mathbf{h}_t$
*input*

$\mathbf{c}_t^{KV}$
*compressed latent vector*

$W^{UK}$

$\{\mathbf{k}_{t,i}^C\}$
*key*

$W^{UV}$

$\{\mathbf{v}_{t,i}^C\}$
*value*

down-project

up-project

*t: token position index*
*i: attention head index*

# Multi-head Latent Attention (MLA)

Do low-rank compression for key K and value V to a single vector and do up-projection to the original vector space during inference.

$$d_c << d_h n_h$$

Down-project to C
(C: head-agnostic)
$$C_t^{KV} = W^{DKV} h_t \qquad \mathbf{c}_t^{KV} \in \mathbb{R}^{d_c} \quad W^{DKV} \in \mathbb{R}^{d_c \times d}$$

Up-project to K
$$K_t^C = W^{UK} C_t^{KV}$$

$$W^{UK}, W^{UV} \in \mathbb{R}^{d_h n_h \times d_c}$$

Up-project to V
$$V_t^C = W^{UV} C_t^{KV}$$

D: down-project
U: up-project

# Multi-head Latent Attention (MLA)

# Multi-head Latent Attention (MLA)

Full Picture:

- Query: In order to reduce the activation memory during training, they also perform low-rank compression for query Q, even if it cannot reduce the KV cache.
- Decoupled RoPE:Then use these additional vectors in the RoPE process only, while keeping the original keys isolated with the rotation matrix, and concatenate in the end.

$$\mathbf{c}_t^Q = W^{DQ}\mathbf{h}_t,$$

$$[\mathbf{q}_{t,1}^C; \mathbf{q}_{t,2}^C; ...; \mathbf{q}_{t,n_h}^C] = \mathbf{q}_t^C = W^{UQ}\mathbf{c}_t^Q,$$

$$[\mathbf{q}_{t,1}^R; \mathbf{q}_{t,2}^R; ...; \mathbf{q}_{t,n_h}^R] = \mathbf{q}_t^R = \text{RoPE}(W^{QR}\mathbf{c}_t^Q),$$

$$\mathbf{q}_{t,i} = [\mathbf{q}_{t,i}^C; \mathbf{q}_{t,i}^R],$$

$$\boxed{\mathbf{c}_t^{KV}} = W^{DKV}\mathbf{h}_t,$$

$$[\mathbf{k}_{t,1}^C; \mathbf{k}_{t,2}^C; ...; \mathbf{k}_{t,n_h}^C] = \mathbf{k}_t^C = W^{UK}\mathbf{c}_t^{KV},$$

$$\boxed{\mathbf{k}_t^R} = \text{RoPE}(W^{KR}\mathbf{h}_t),$$

$$\mathbf{k}_{t,i} = [\mathbf{k}_{t,i}^C; \mathbf{k}_t^R],$$

$$\mathbf{o}_{t,i} = \sum_{j=1}^{t} \text{Softmax}_j\left(\frac{\mathbf{q}_{t,i}^T\mathbf{k}_{j,i}}{\sqrt{d_h + d_h^R}}\right)\mathbf{v}_{j,i}^C,$$

$$\mathbf{u}_t = W^O[\mathbf{o}_{t,1}; \mathbf{o}_{t,2}; ...; \mathbf{o}_{t,n_h}],$$

blue: cached components

# Multi-head Latent Attention (MLA)

Performance vs classic MHA
- Better memory efficiency
- Better model capacity

| Attention Mechanism | KV Cache per Token (# Element) | Capability |
|---|---|---|
| Multi-Head Attention (MHA) | $2n_h d_h l$ | Strong |
| Grouped-Query Attention (GQA) | $2n_g d_h l$ | Moderate |
| Multi-Query Attention (MQA) | $2d_h l$ | Weak |
| MLA (Ours) | $(d_c + d_h^R)l \approx \frac{9}{2}d_h l$ | Stronger |

# Multi-head Latent Attention (MLA)

Performance vs classic MHA
- Better memory efficiency
- Better model capacity

| Benchmark (Metric) | # Shots | Small MoE w/ MHA | Small MoE w/ MLA | Large MoE w/ MHA | Large MoE w/ MLA |
|---|---|---|---|---|---|
| # Activated Params | - | 2.5B | 2.4B | 25.0B | 21.5B |
| # Total Params | - | 15.8B | 15.7B | 250.8B | 247.4B |
| KV Cache per Token (# Element) | - | 110.6K | 15.6K | 860.2K | 34.6K |
| BBH (EM) | 3-shot | 37.9 | **39.0** | 46.6 | **50.7** |
| MMLU (Acc.) | 5-shot | 48.7 | **50.0** | 57.5 | **59.0** |
| C-Eval (Acc.) | 5-shot | **51.6** | 50.9 | 57.9 | **59.2** |
| CMMLU (Acc.) | 5-shot | 52.3 | **53.4** | 60.7 | **62.5** |

# DeepSeek MoE



**Text Prediction** | **Task Classifier**

**5. Output Representation**

Layer Norm

**3. MLP** → Feed Forward

12x

Layer Norm

**2. Self-Attention** → Masked Multi Self Attention

**4. Repeat**

**1. Input Embeddings (each token = vector)** → Text & Position Embed

Decoder-only transformer (GPT)

# DeepSeek MoE: Feed Forward



Most of the parameters are here

Decoder-only transformer (GPT)

# DeepSeek MoE: Architecture

- Imagine the human brain.
- We have distinct parts responsible for different tasks.
- For a specific task, we can invoke just the specific part of the brain to do the work
- This approach is more scalable & efficient since you only use a subset of available neurons
- Instead of having one FFN, we have multiple (relatively) smaller FFNs (now called experts)

# DeepSeek MoE: Architecture

- Previous MoE architectures had limited experts (N=2 to 16) with only ~2 of them being active for a token
- This did bring the efficiency gains.
- But, there are two problems:
  - Knowledge Hybridity:
    - Each expert is not really specific to a certain task.
    - It's still quite big and holds diverse, hybrid knowledge.
  - Knowledge Redundancy
    - Same information is stored across multiple experts.
- **2-16 experts is not enough!**

# DeepSeek MoE: Architecture

- To combat this, Deepseek introduces two new innovations:
  - Fine-grained expert segmentation (combat K. Hybridity)
    - Each expert is split into multiple smaller ones, keeping parameter count the same.
    - With more experts, knowledge can be decomposed and each expert can get more specialized in a certain area.
  - Shared experts
    - Some experts are always activated and can focus on common knowledge (combat K. Redundancy)

# DeepSeek MoE

# DeepSeek MoE

Output Hidden $\mathbf{h}'_t$

Routed Expert

Shared Expert

A learned function that provides a weight to remaining experts

Router

Top-$K_r$

Input Hidden $\mathbf{u}_t$

Output Hidden $\mathbf{h}'_t$

Routed Expert

Shared Expert

Router

Top-$K_r$

Input Hidden $\mathbf{u}_t$

Input is passed to shared + chosen experts

# DeepSeek MoE



Output from each routed expert is weighted by the weight assigned to it by the router

Output Hidden $\mathbf{h}_t'$

Combine output from all the experts

$1$ ... $N_s$

$1$ $2$ $3$ $4$ ... $N_r$-1 $N_r$

Router

Top-$K_r$

Input Hidden $\mathbf{u}_t$

Output Hidden $\mathbf{h}'_t$

Input for the next transformer block...

$1$ ... $N_s$   $1$   $2$   $3$   $4$ ... $N_r$-1   $N_r$

Router   Top-$K_r$

Input Hidden $\mathbf{u}_t$

# DeepSeek MoE

$\mathbf{u_t}$   FFN-Input of the t-th token (output from self-attention layer)

$\mathbf{e_i}$   Centroid vector of the i-th routed expert

(you can imagine this as a learned vector representation of the kind of tokens a certain expert is meant to deal with)

Taking a Dot Product

(How similar is this token to the kind of tokens expert-i is used to dealing with?)

$$\mathbf{u}_t^T \mathbf{e}_i$$

# DeepSeek MoE

$\mathbf{u_t}$    FFN-Input of the t-th token (output from self-attention layer)

$\mathbf{e_i}$    Centroid vector of the i-th routed expert

<span style="color:red">Pass the similarity (affinity) score of token 't' with expert 'i' through sigmoid to get a value from 0 to 1</span>

$$s_{i,t} = \text{Sigmoid}\left(\mathbf{u}_t{}^T \mathbf{e}_i\right)$$

# DeepSeek MoE

$\mathbf{u_t}$    FFN-Input of the t-th token (output from self-attention layer)

$\mathbf{e_i}$    Centroid vector of the i-th routed expert

$s_{i,t}$    Affinity score of token-t with expert-i

<span style="color:red">Choose only TopK (K_r) experts based on affinity scores</span>

$$
g'_{i,t} = \begin{cases} s_{i,t}, & s_{i,t} \in \text{Topk}(\{s_{j,t} | 1 \leqslant j \leqslant N_r\}, K_r), \\ 0, & \text{otherwise,} \end{cases}
$$

# DeepSeek MoE

$\mathbf{u_t}$    FFN-Input of the t-th token (output from self-attention layer)

$\mathbf{e_i}$    Centroid vector of the i-th routed expert

$s_{i,t}$    Affinity score of token-t with expert-i

$$g_{i,t} = \frac{g'_{i,t}}{\sum_{j=1}^{N_r} g'_{j,t}}$$

Run softmax (without the exponential) over the N_r routed experts, getting a probability distribution over the routed experts. This is referred to as the 'gating value'

# DeepSeek MoE

$\mathbf{u_t}$  FFN-Input of the t-th token (output from self-attention layer)

$g_{i,t}$  Gating value for expert-i with token-t

Putting everything together...

# DeepSeek MoE

$\mathbf{u_t}$    FFN-Input of the t-th token (output from self-attention layer)

$g_{i,t}$    Gating value for expert-i with token-t

$$\mathbf{h}'_t = \mathbf{u}_t + \sum_{i=1}^{N_s} \text{FFN}_i^{(s)}(\mathbf{u}_t) + \sum_{i=1}^{N_r} g_{i,t} \text{FFN}_i^{(r)}(\mathbf{u}_t)$$

Pass input through the N_s shared experts and combine them

# DeepSeek MoE

$\mathbf{u_t}$ FFN-Input of the t-th token (output from self-attention layer)

$g_{i,t}$ Gating value for expert-i with token-t

$$\mathbf{h}'_t = \mathbf{u}_t + \sum_{i=1}^{N_s} \text{FFN}_i^{(s)}(\mathbf{u}_t) + \sum_{i=1}^{N_r} g_{i,t} \text{FFN}_i^{(r)}(\mathbf{u}_t)$$
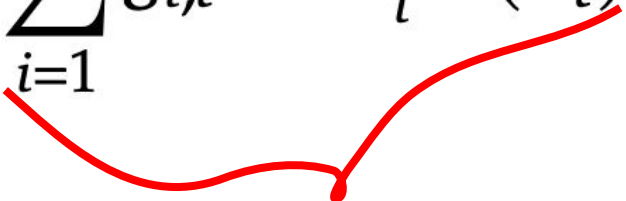
Pass input through the N_r routed experts and do weighted sum based on gating value

# DeepSeek MoE

$\mathbf{u_t}$    FFN-Input of the t-th token (output from self-attention layer)

$g_{i,t}$    Gating value for expert-i with token-t

$$\mathbf{h}'_t = \mathbf{u}_t + \sum_{i=1}^{N_s} \text{FFN}_i^{(s)}(\mathbf{u}_t) + \sum_{i=1}^{N_r} g_{i,t} \text{FFN}_i^{(r)}(\mathbf{u}_t)$$

Add the residual (original input)

# DeepSeek MoE

- DeepSeek v3 has 61 transformer blocks.
- Except first 3 FFNs, all are substituted with MoE.
- Each MoE layer has **1 shared expert** and **256 routed experts.**
  - Some details were abstracted earlier. In actual, we have 8 groups (large experts) among which 2 are chosen first. So, each group has 256 / 8 = 32 smaller experts.
- Among the routed experts, **8** will be chosen for each token.
- Intermediate dimension of each expert is just **2048.**
- So, although, DeepSeek V3 comprises of 671B parameters, **only 37B are activated for each token.**

# Conclusion

- **Multi-Head Latent Attention (MLA)** dramatically reduces memory usage compared to traditional attention mechanisms
- **Mixture of Experts (MoE)** architecture activates only 37B of 671B total parameters per token, optimizing resource utilization
- DeepSeek v3 achieves state-of-the-art performance while significantly lowering computational costs

# References

1.  DeepSeek-AI. (2024). DeepSeek-V3 Technical Report. ([link](link))
2.  DeepSeek-AI. (2024). DeepSeek-V2: A Strong, Economical, and Efficient Mixture-of-Experts Language Model. ([link](link))
3.  Dai, D., Deng, C., Zhao, C., Xu, R. X., & Gao, H. (2024). DeepSeekMoE: Towards Ultimate Expert Specialization in Mixture-of-Experts Language Models. ([link](link))
4.  Jiang, A. Q., Sablayrolles, A., Roux, A., Mensch, A., & Savary, B. (2024). Mixtral of Experts ([link](link)).