



香港科技大學
THE HONG KONG
UNIVERSITY OF SCIENCE
AND TECHNOLOGY

COMP 4901B

Large Language Models

MoE LLMs and Review

Junxian He

Nov 26, 2025

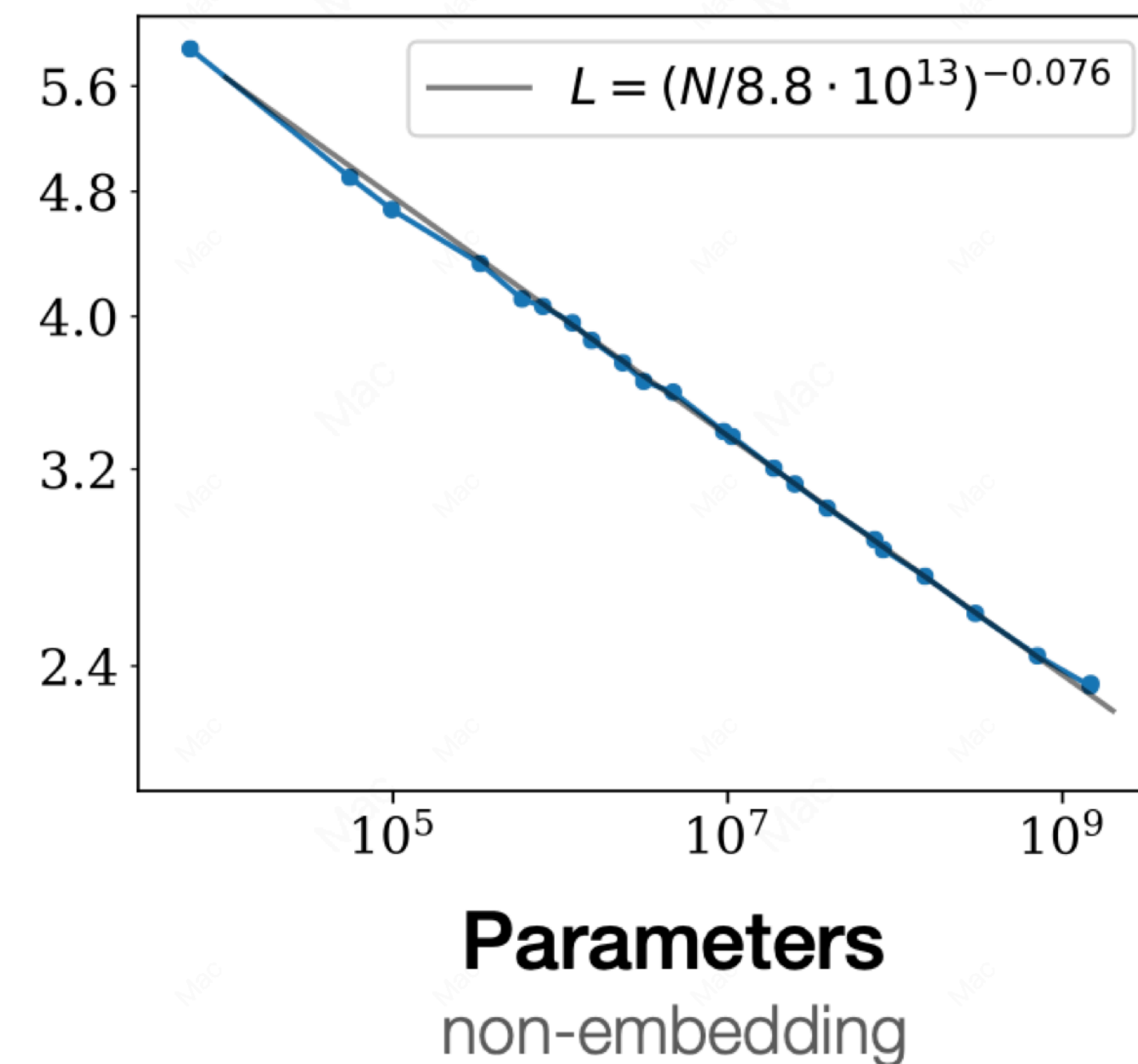
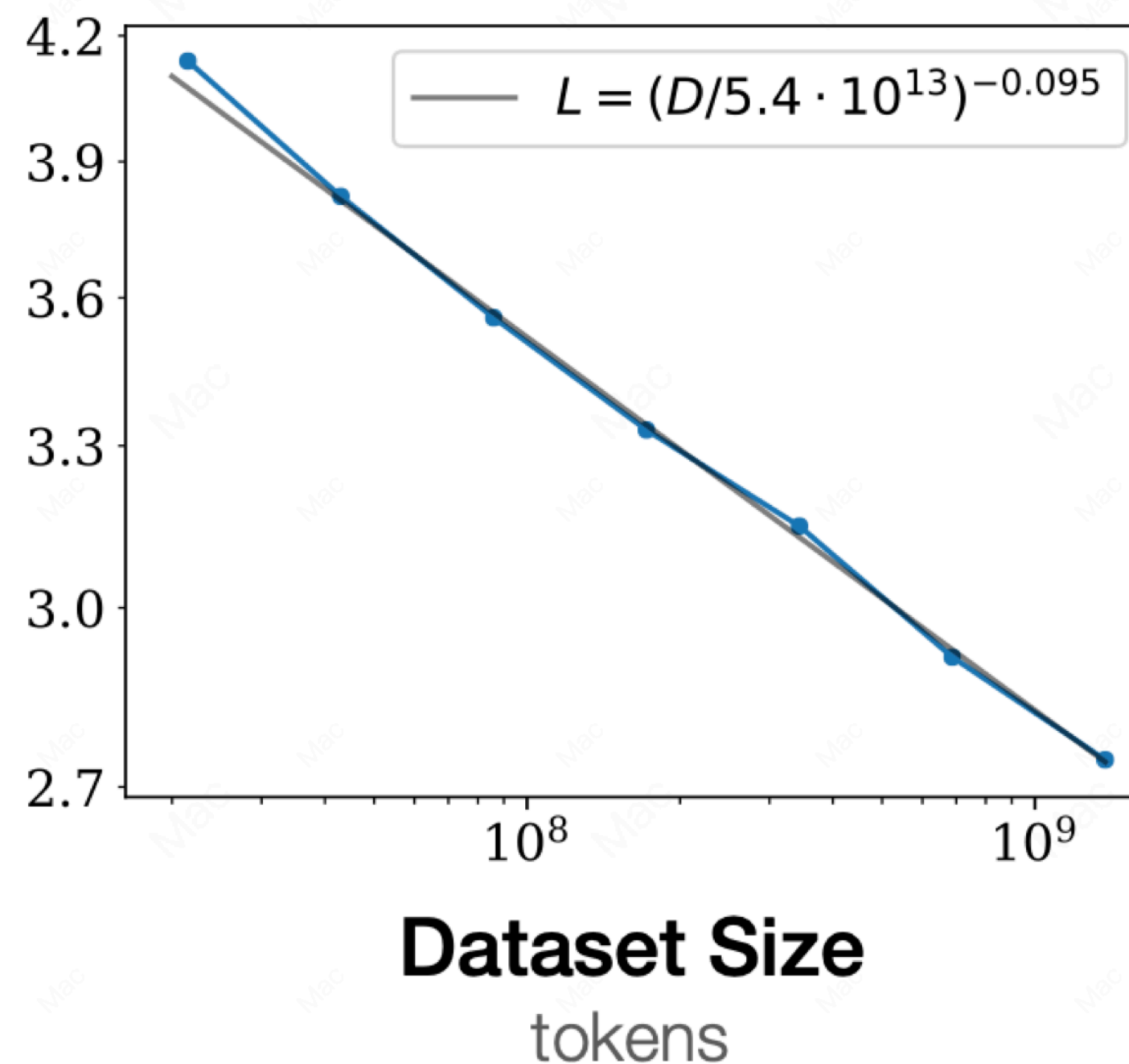
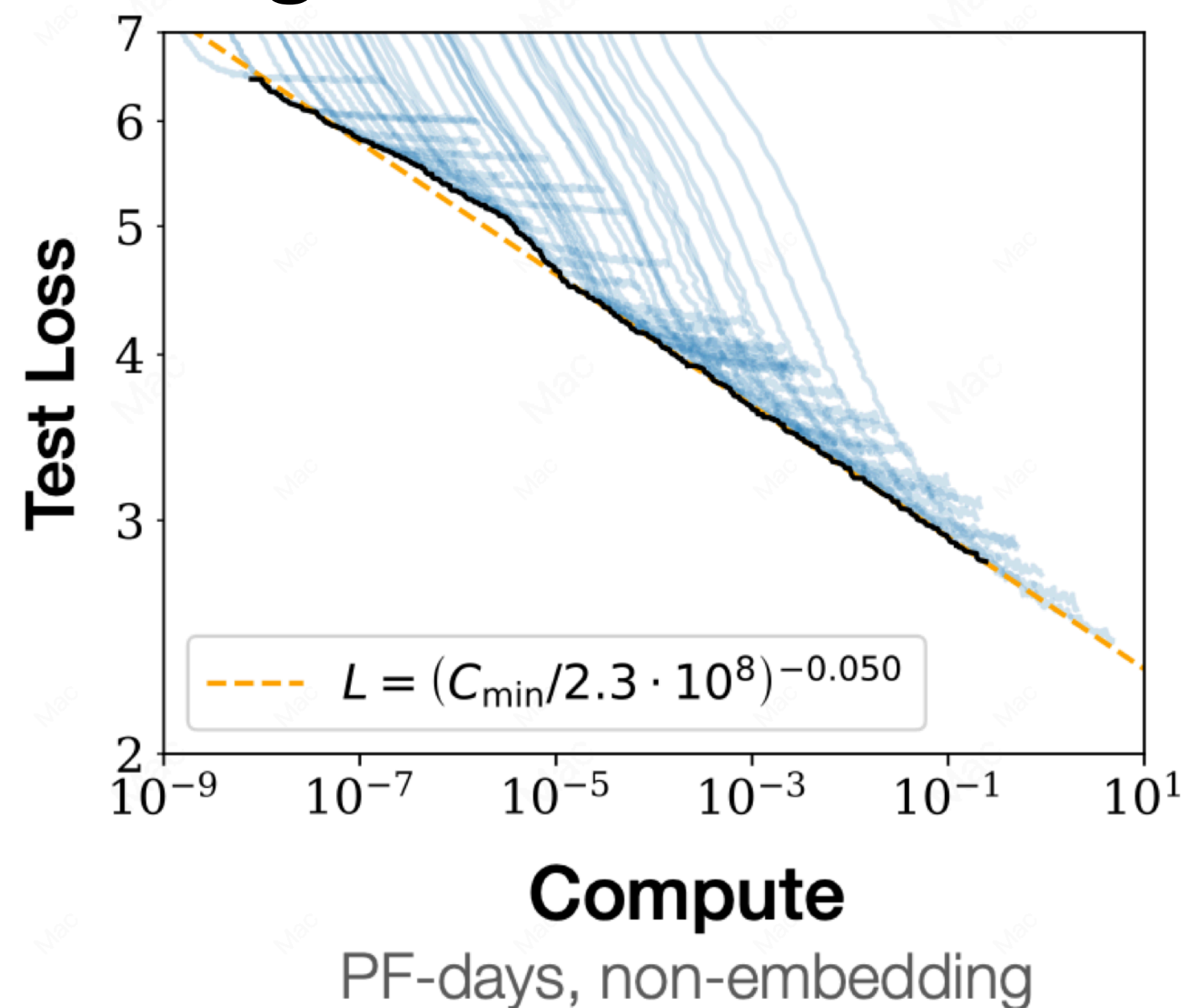
Final exam:

Dec 11 (Thurs) 8:30am-10:30am

LTC

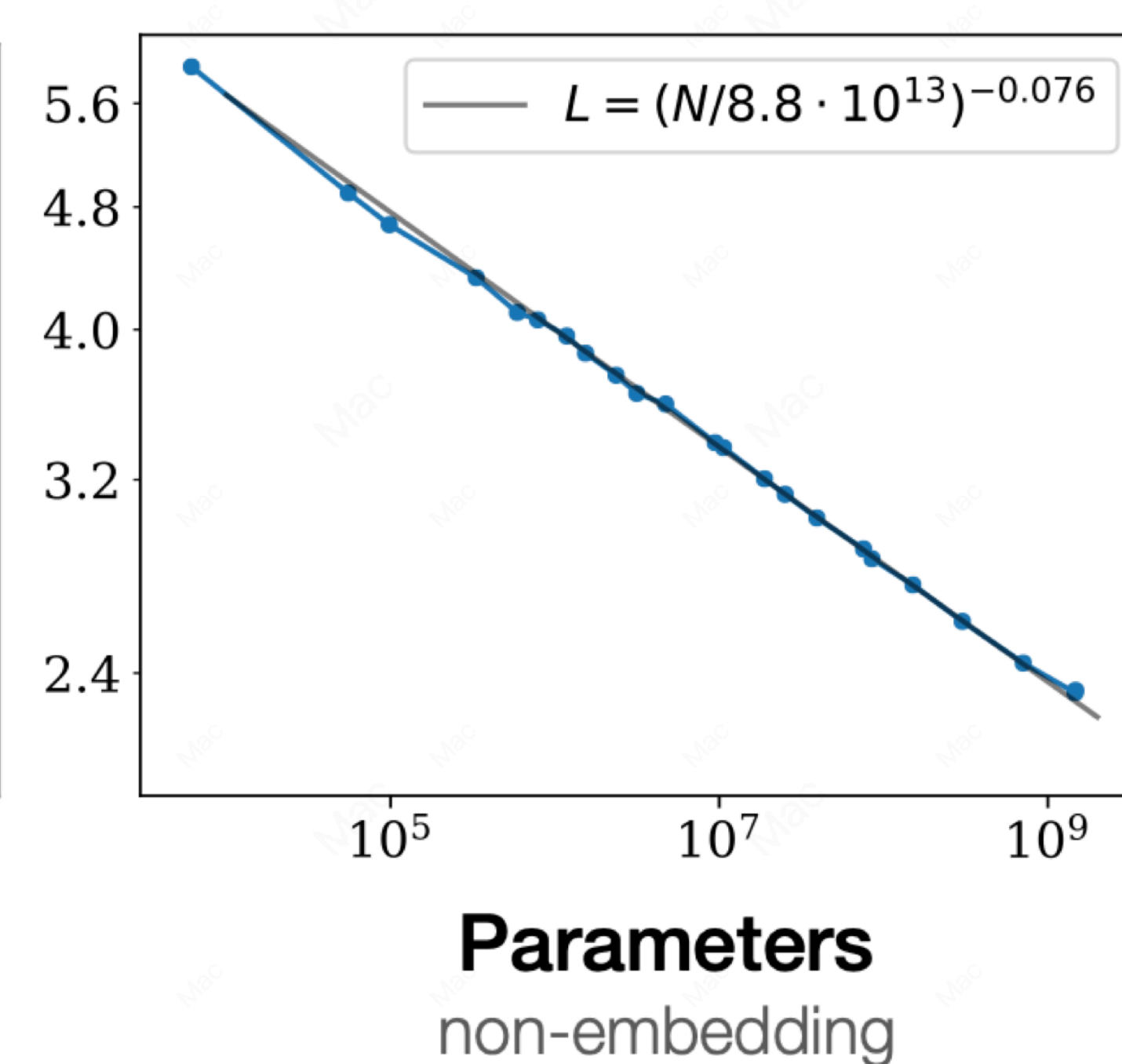
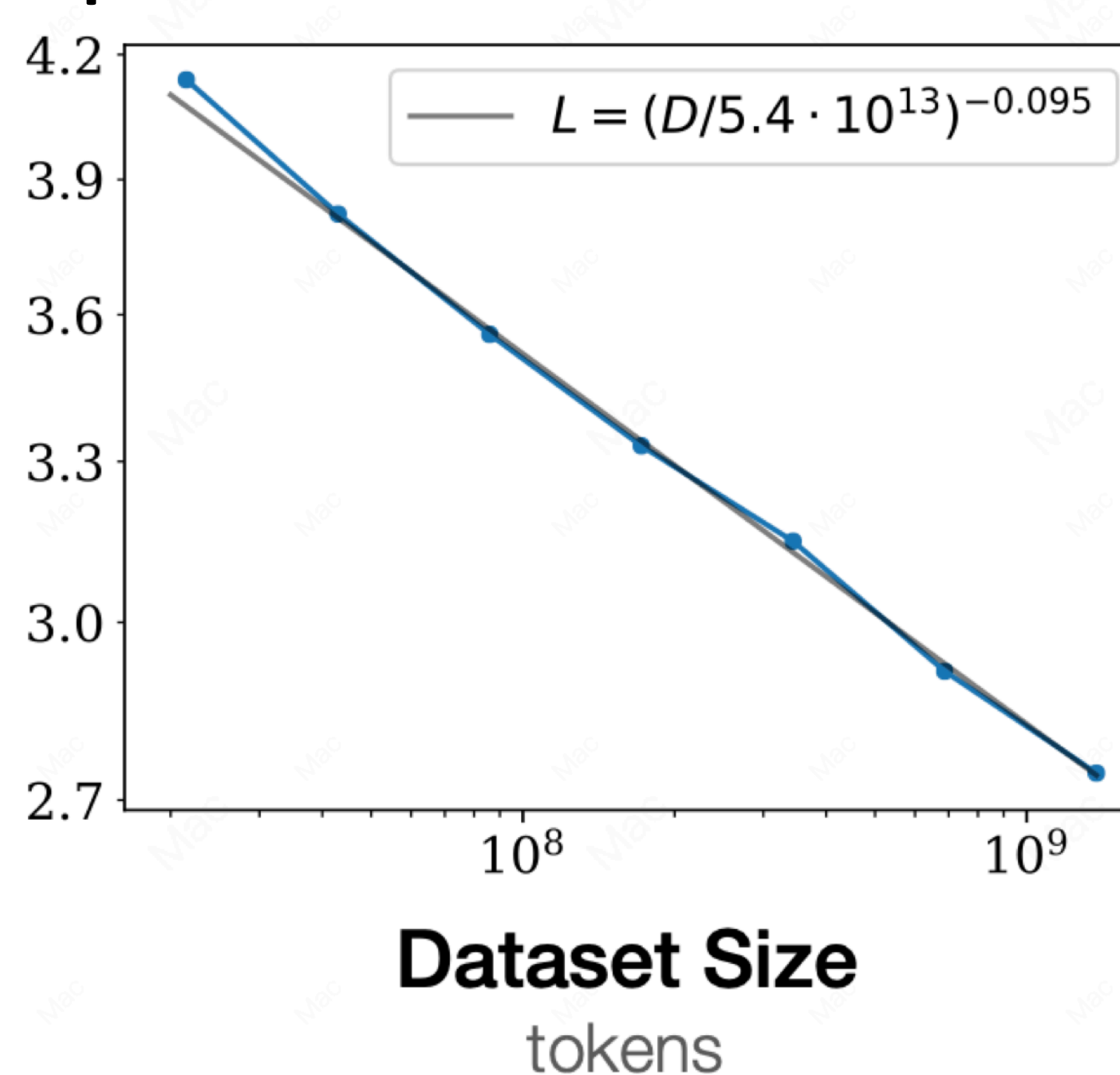
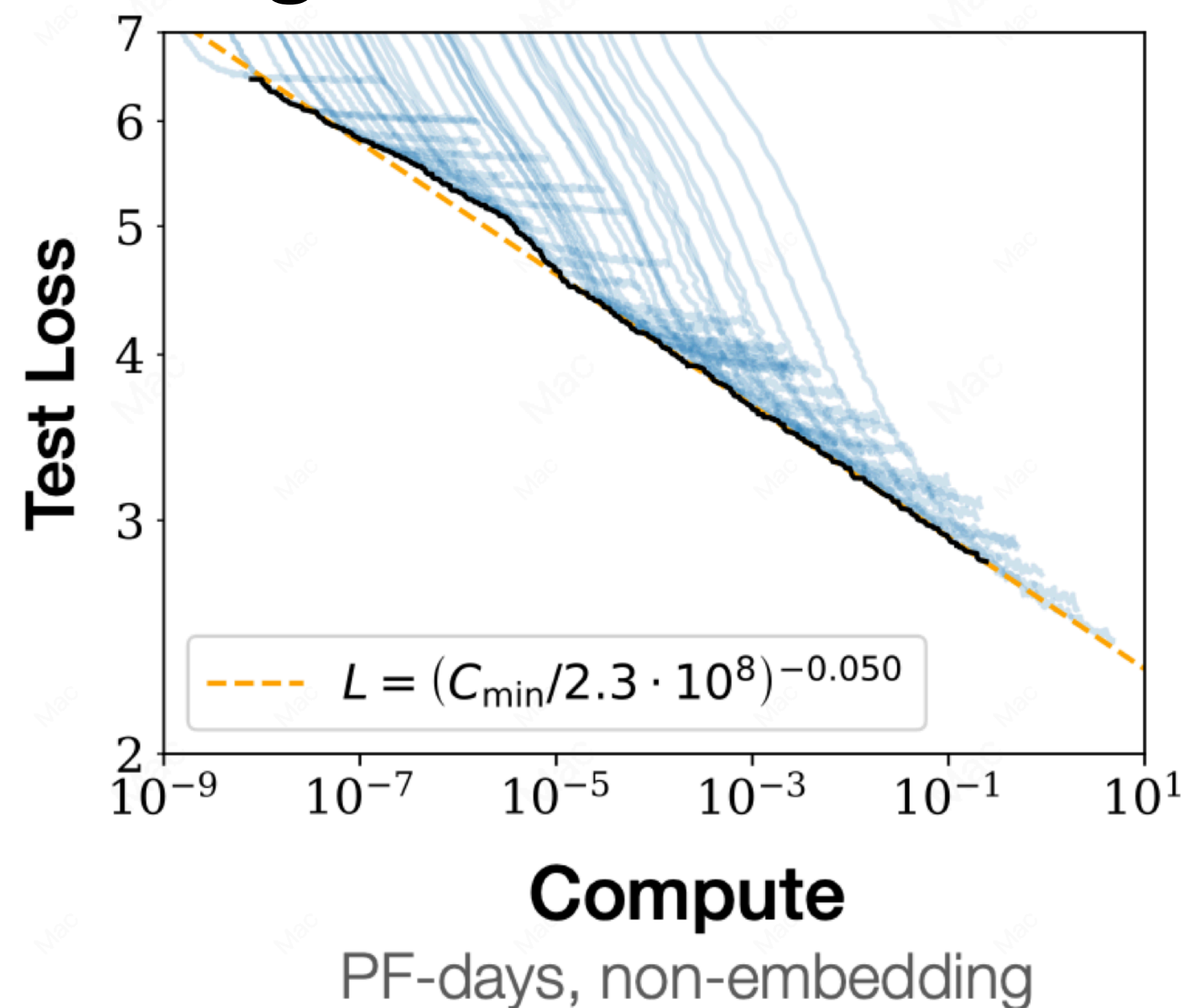
Challenges of Scaling Model Sizes

Scaling law tells us to scale up model sizes



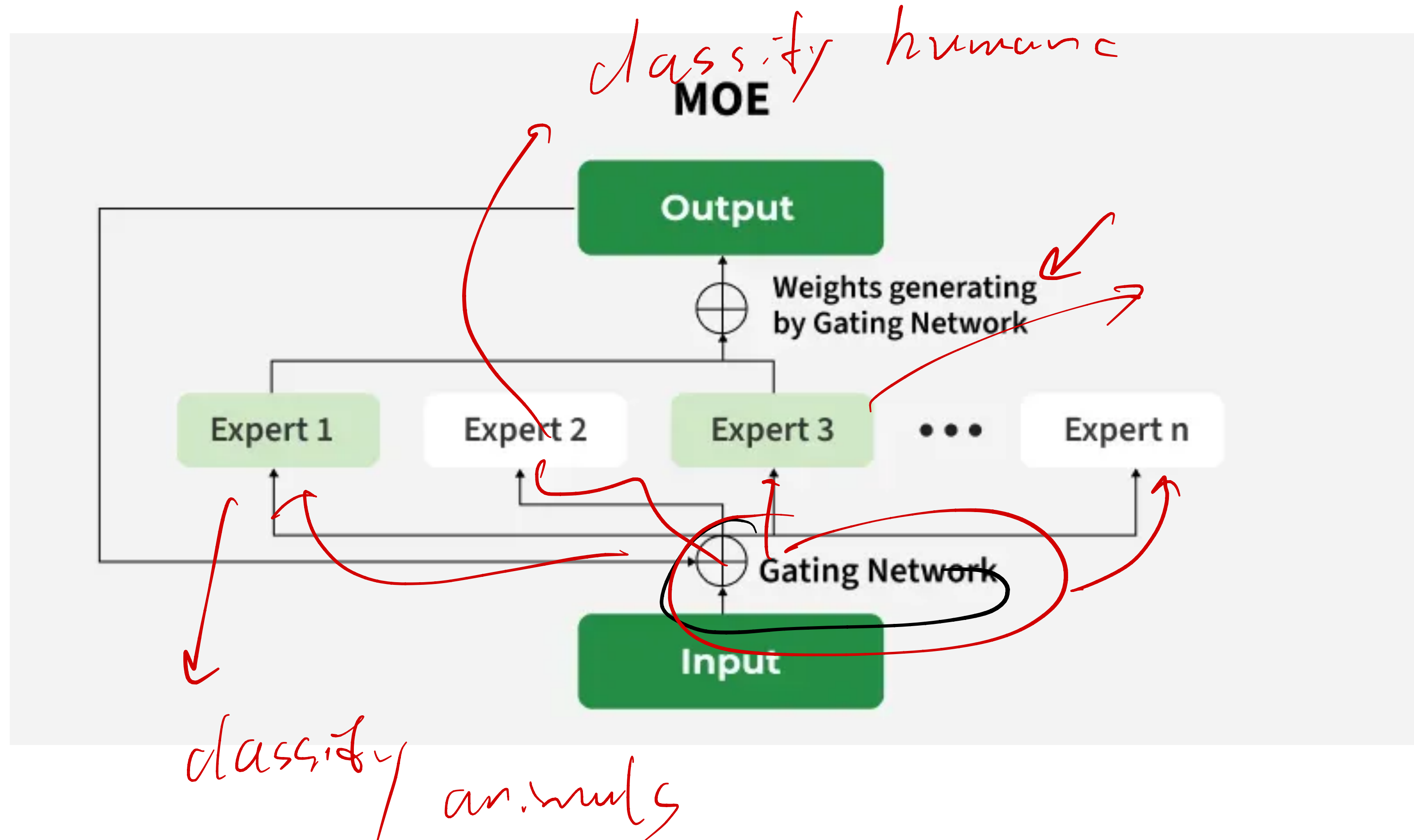
Challenges of Scaling Model Sizes

Scaling law tells us to scale up model sizes

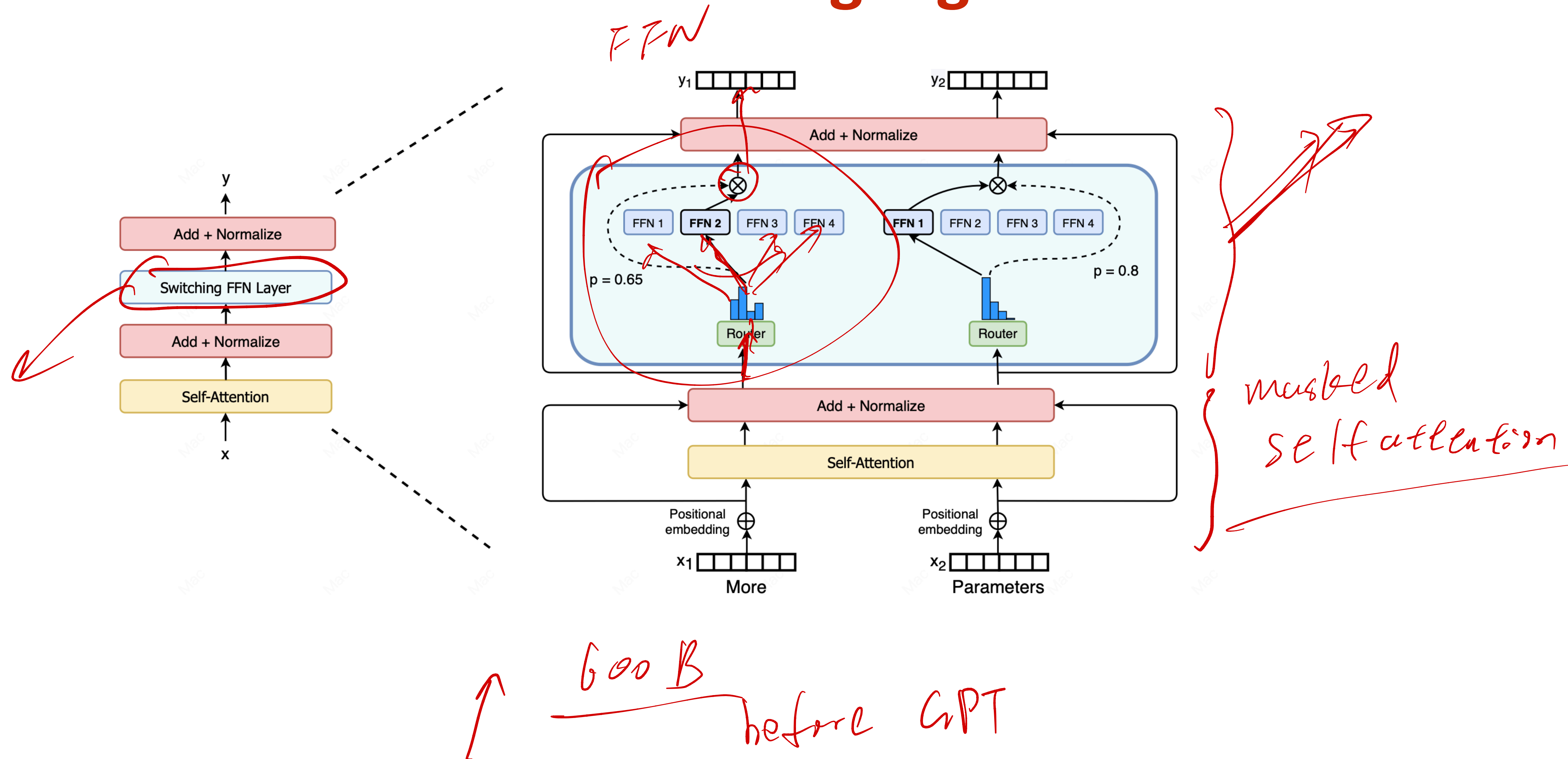


However, larger model sizes require more compute to train and causes higher latency

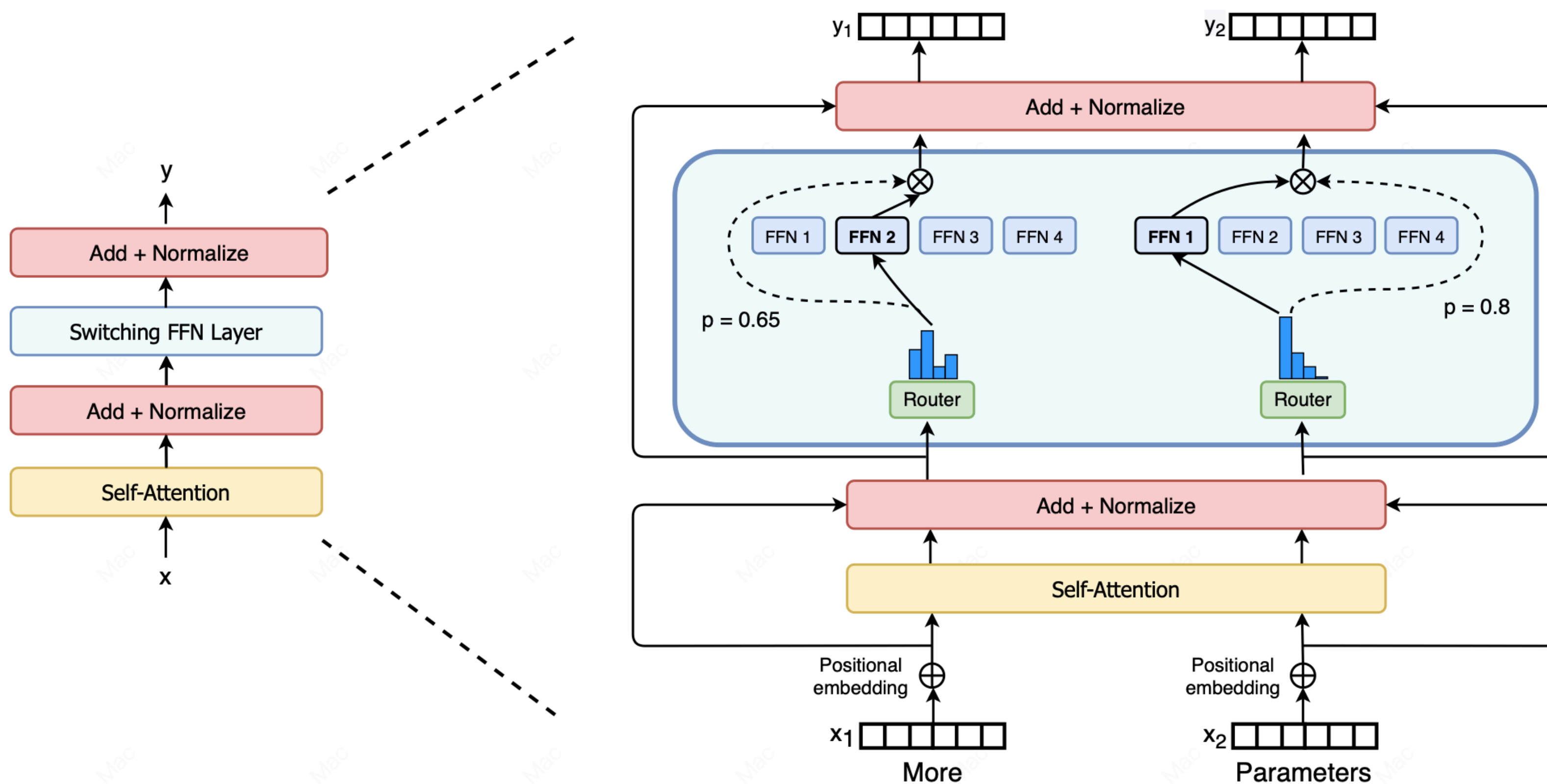
Mixture of Experts (MoE) in Traditional Machine Learning



MoE Transformer Language Models

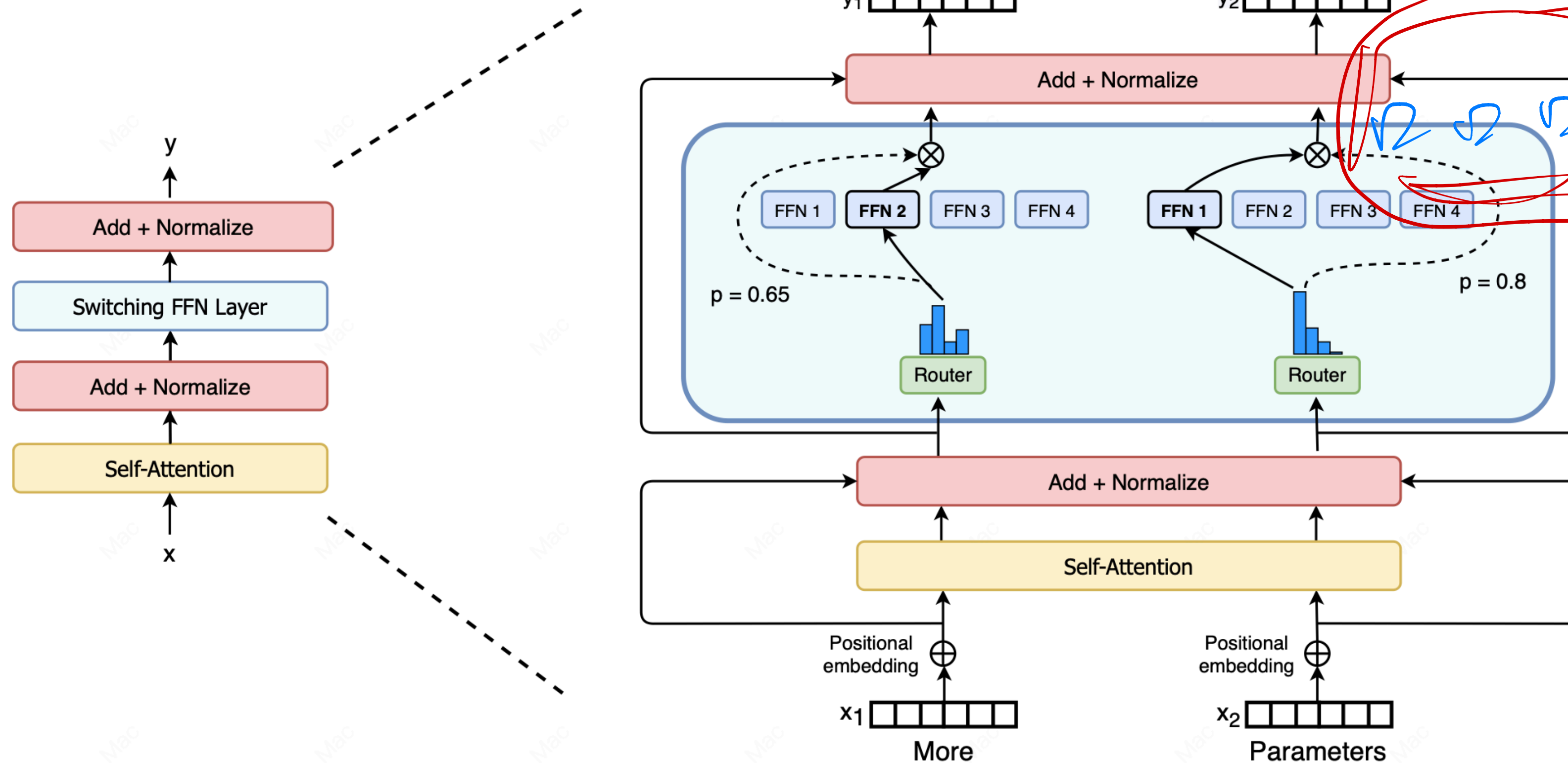


MoE Transformer Language Models



Mixture of FFN Blocks

MoE Transformer Language Models



Handwritten notes:

- Experts** (circled in blue)
- Router** (circled in blue)
- 2 or 3** (circled in red)

Mixture of FFN Blocks

For each token at each layer, only a small fraction (e.g., 2 or 3) experts are activated by the router, thus this is also referred to as SPARSE models

Sparse Routing

Sparse Routing

parameter
router

$$h(x) = W_r \cdot x$$

Logits of different experts

Sparse Routing

$$\underline{h(x)} = W_r \cdot x \quad \text{Logits of different experts}$$

$p_i(x)$ prob x
routing to
expert i

$$p_i(x) = \frac{e^{h(x)_i}}{\sum_j^N e^{h(x)_j}} \quad \text{Gate value, this is softmax}$$

Sparse Routing

$$h(x) = W_r \cdot x \quad \text{Logits of different experts}$$

$$p_i(x) = \frac{e^{h(x)_i}}{\sum_j^N e^{h(x)_j}} \cdot \quad \text{Gate value, this is softmax}$$

Sparse Routing: only top-k experts are used during both training and test time

Sparse Routing

Dense:

$$Ecx$$

$$h(x) = W_r \cdot x \quad \text{Logits of different experts}$$

$$p_i(x) = \frac{e^{h(x)_i}}{\sum_j^N e^{h(x)_j}} \quad \text{Gate value, this is softmax}$$

Sparse Routing: only top-k experts are used during both training and test time

$$y = \sum_{i \in T} p_i(x) E_i(x).$$

aggregation

linear weighted addn

$T = \text{top } k$

Sparse Routing

$$h(x) = W_r \cdot x \quad \text{Logits of different experts}$$

$$p_i(x) = \frac{e^{h(x)_i}}{\sum_j^N e^{h(x)_j}}. \quad \text{Gate value, this is softmax}$$

Sparse Routing: only top-k experts are used during both training and test time

$$y = \sum_{i \in \mathcal{T}} p_i(x) E_i(x).$$

Weighted addition of each expert's output

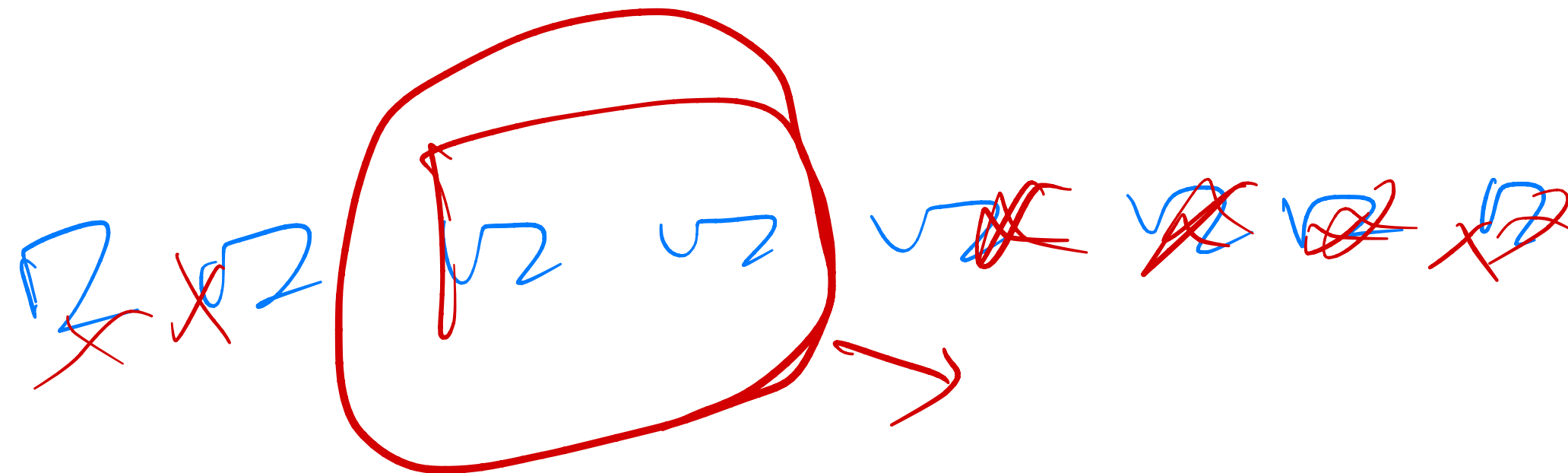
Load Balancing

Only next-token prediction loss may learn to only use a few experts

Load Balancing

Only next-token prediction loss may learn to only use a few experts

Vicious cycle: When a certain subset experts are chosen, next-token prediction will optimize them to suit the input, then these experts are more likely to be chosen



Auxiliary Load Balancing Loss

$$\text{loss} = \alpha \cdot N \cdot \sum_{i=1}^N f_i \cdot P_i$$

where f_i is the fraction of tokens dispatched to expert i ,

$$f_i = \frac{1}{T} \sum_{x \in \mathcal{B}} \mathbb{1}\{\text{argmax } p(x) = i\}$$

and P_i is the fraction of the router probability allocated for expert i ,²

$$P_i = \frac{1}{T} \sum_{x \in \mathcal{B}} p_i(x).$$

$\mathcal{B} = \text{batch}$

learnable

f_i

P_i

optimize

minimize

ideal:

specific routing each token

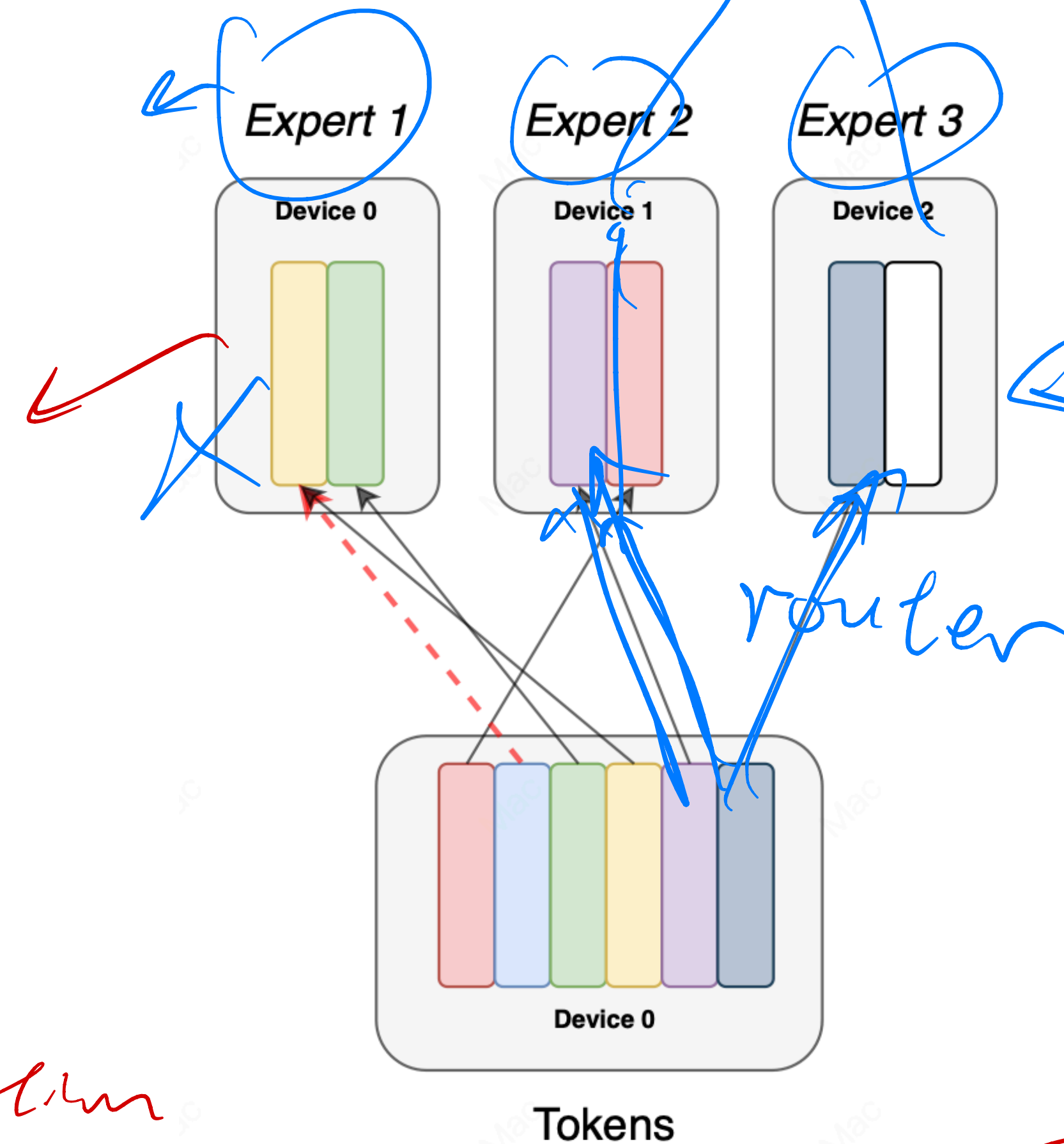
uniform load globally

Why MoE?

not interpretable

aggre

MoE models support easier model parallel across different GPUs. It can easily split models



system infra

router

FFN

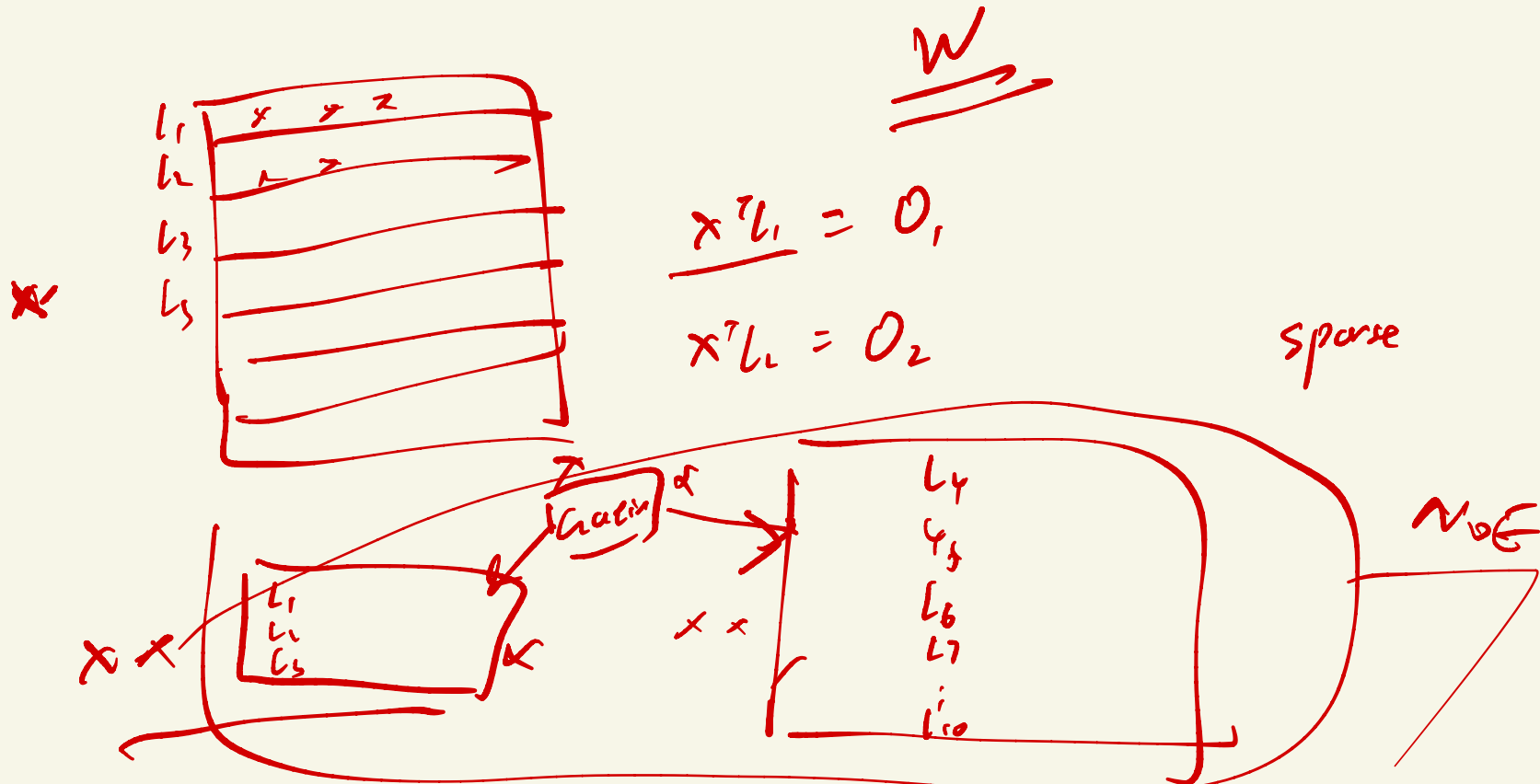
w

activation

w

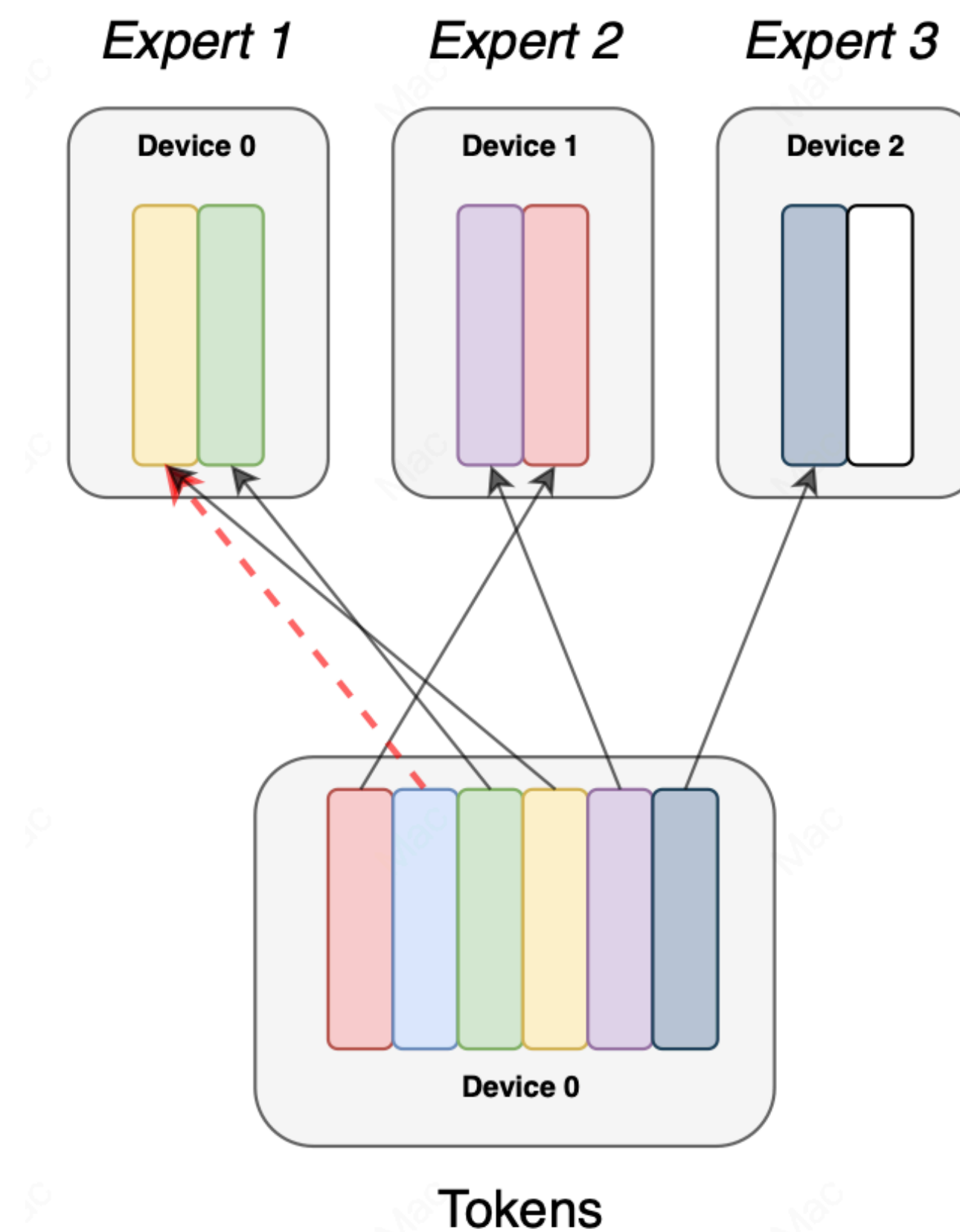
Tokens

Matrix multiplication



Why MoE?

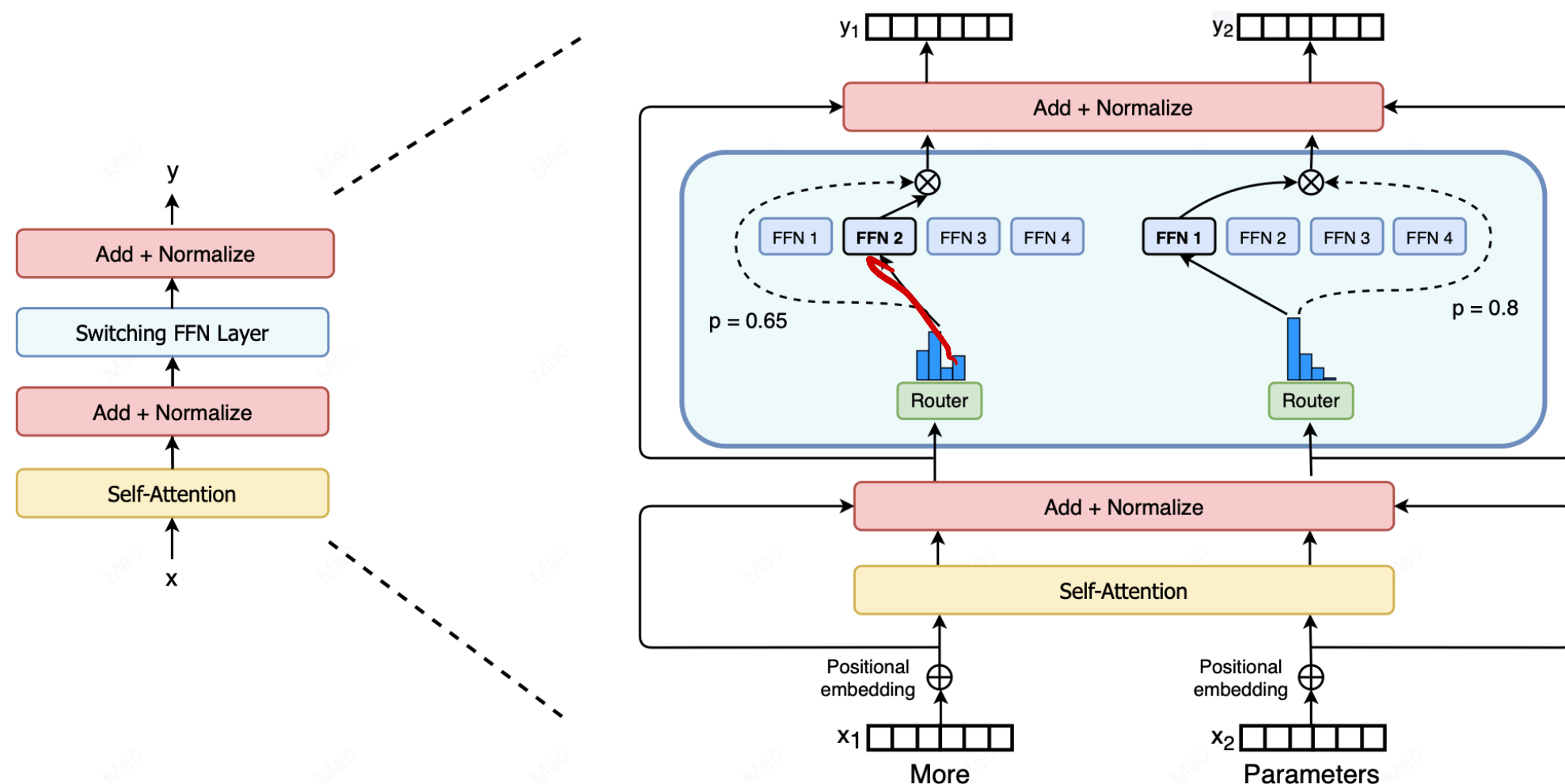
MoE models support easier model parallel across different GPUs. It can easily split models



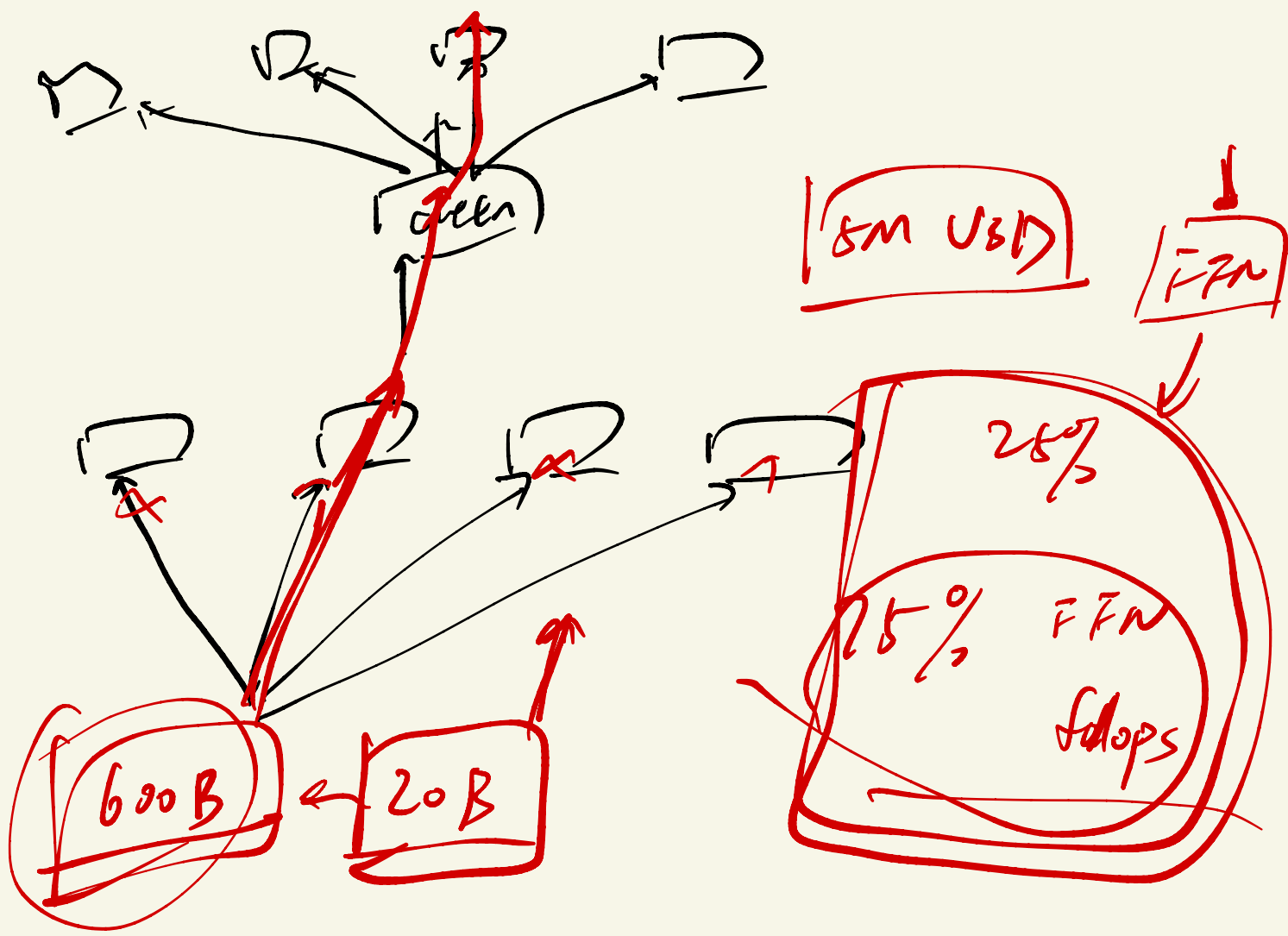
MoE models have major infra-wise benefits when scaling compared to dense models

Why MoE?

forward backmark

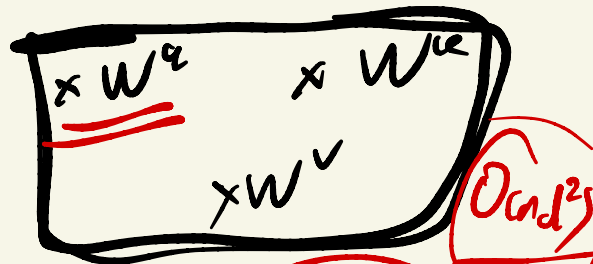


1. Training flops compared to dense model of the same size?
2. Inference flops compared to dense model of the same size?



atten

d : hidden size
 n : seq length.



QK

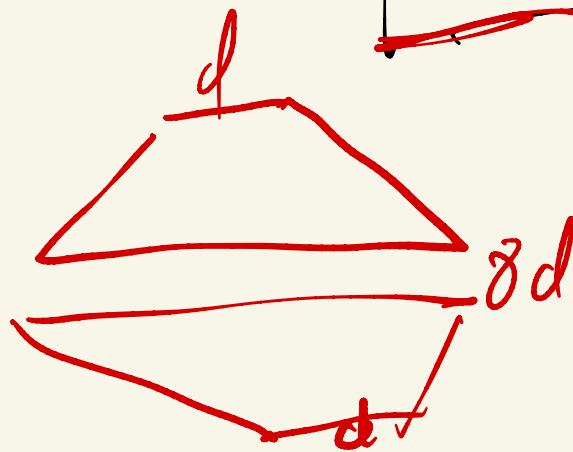
$O(n^2 d)$

FFN

dominant parameter

dominant flops

$$n^2 d = 16 n d^2$$



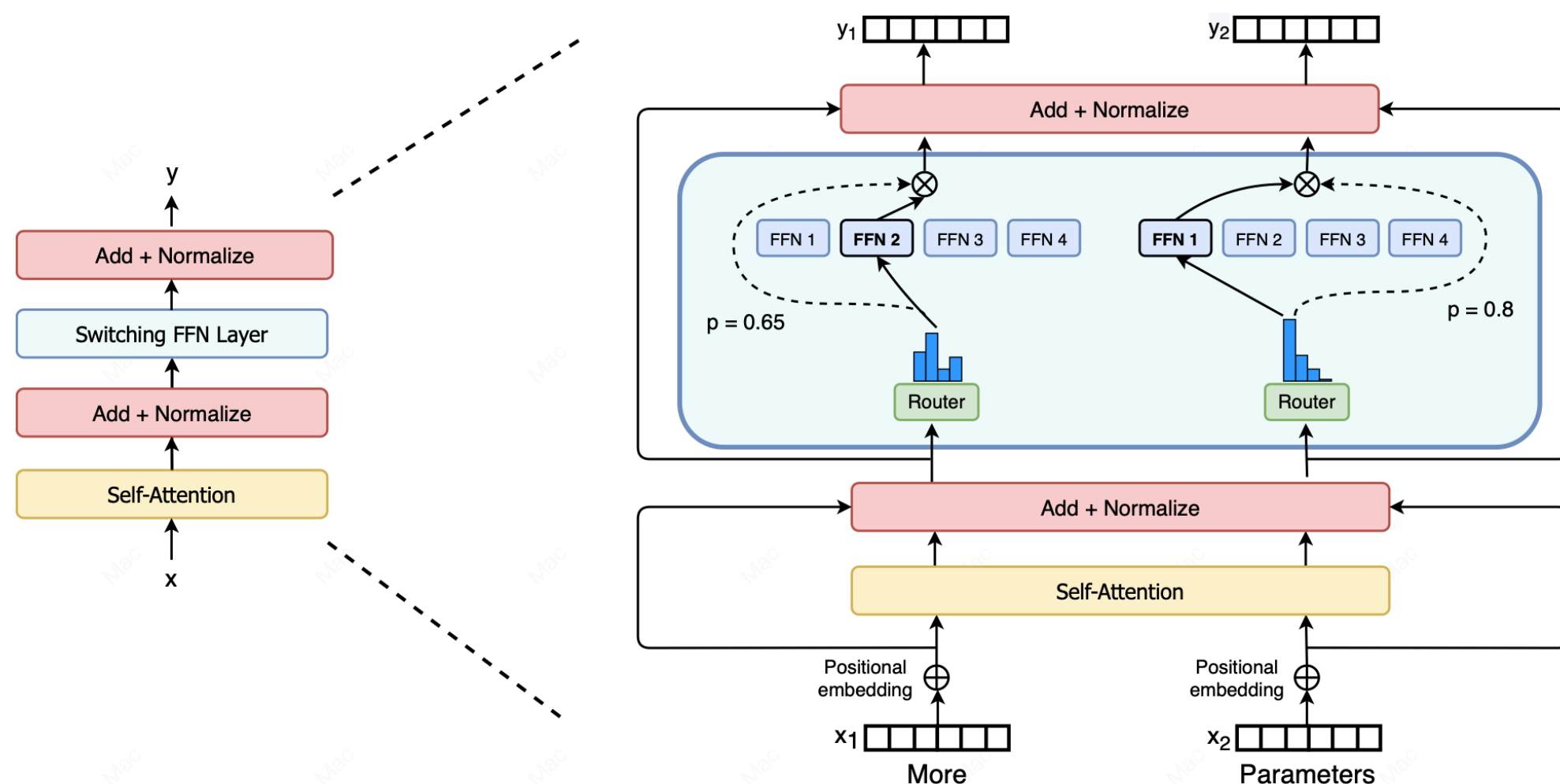
$$O(8 d^2 \times 2 \times n)$$

$$+ 16 n d^2$$

$d \approx 2048$

$n \rightarrow \infty$

Why MoE?

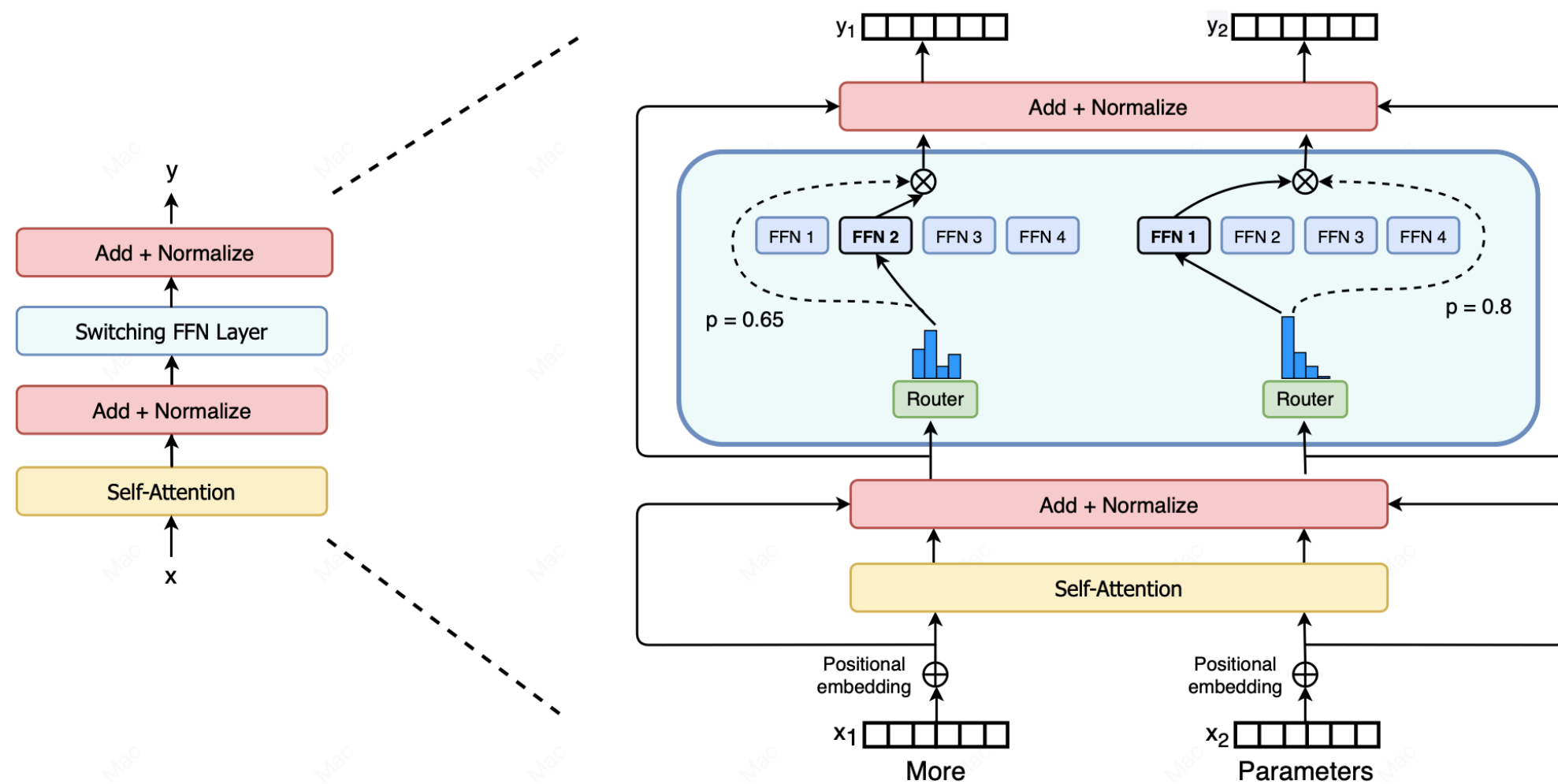


1. Training flops compared to dense model of the same size?
2. Inference flops compared to dense model of the same size?

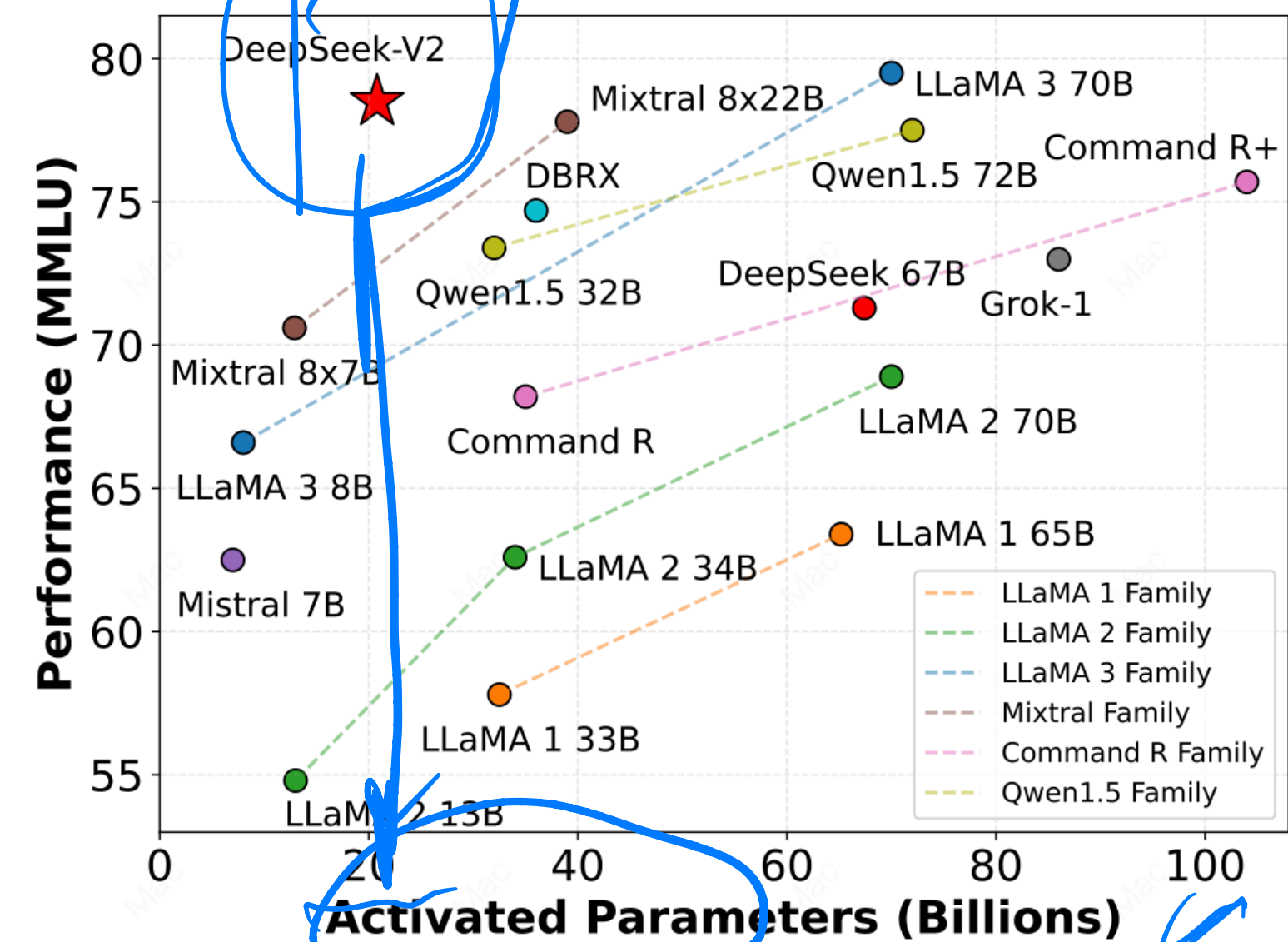
Theoretically, if the activated parameters are only 1B, then the inference latency can be optimized to be close to an 1B dense model

Why MoE?

1. Training flops compared to dense model of the same size?
2. Inference flops compared to dense model of the same size?

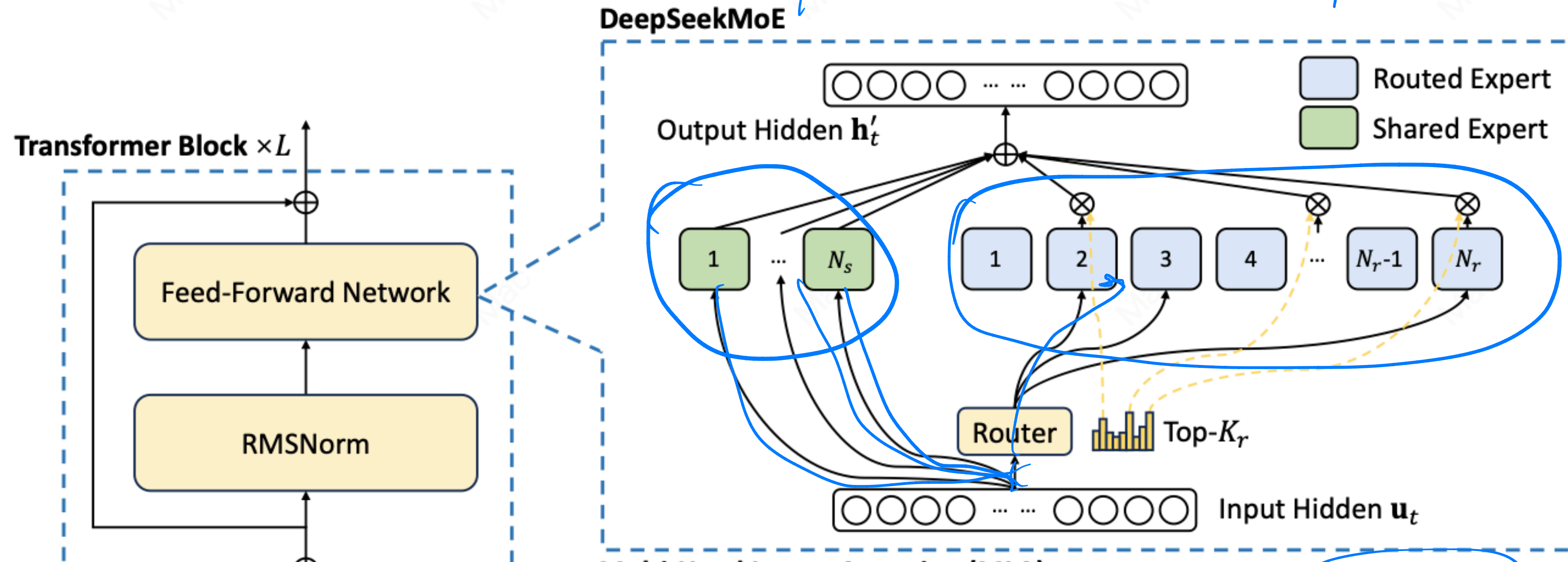


Theoretically, if the activated parameters are only 1B, then the inference latency can be optimized to be close to an 1B dense model



DeepSeek MoE

Deepseek V3 \rightarrow post training



Deepseek R_1

Shared Experts + Routed Expert

For DS-V3, 1 shared expert + 256 routed experts, each token 8 experts are activated

8
256

residual

DeepSeek-v3

Shared experts

Routed experts

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

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

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

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

softmax

input

\mathbf{e}_i router parameter

Loss-Free Load Balancing in DeepSeek-v3

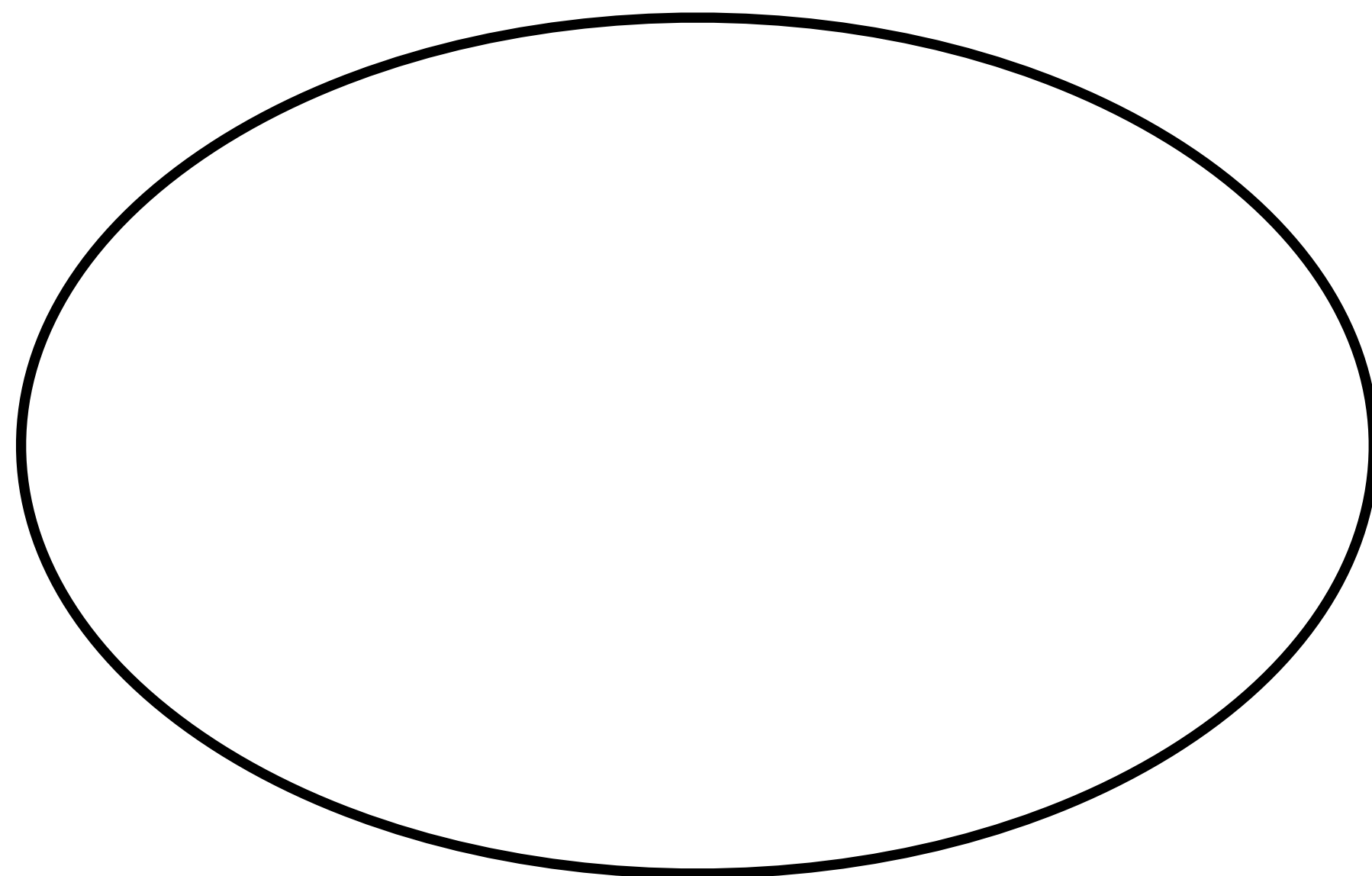
$$g'_{i,t} = \begin{cases} s_{i,t}, & s_{i,t} + b_i \in \text{Topk}(\{s_{j,t} + b_j | 1 \leq j \leq N_r\}, K_r), \\ 0, & \text{otherwise.} \end{cases}$$

Note that the bias term is only used for routing. The gating value, which will be multiplied with the FFN output, is still derived from the original affinity score $s_{i,t}$. During training, we keep monitoring the expert load on the whole batch of each training step. At the end of each step, we will decrease the bias term by γ if its corresponding expert is overloaded, and increase it by γ if its corresponding expert is underloaded, where γ is a hyper-parameter called bias update speed. Through the dynamic adjustment, DeepSeek-V3 keeps balanced expert load during training, and achieves better performance than models that encourage load balance through pure auxiliary losses.

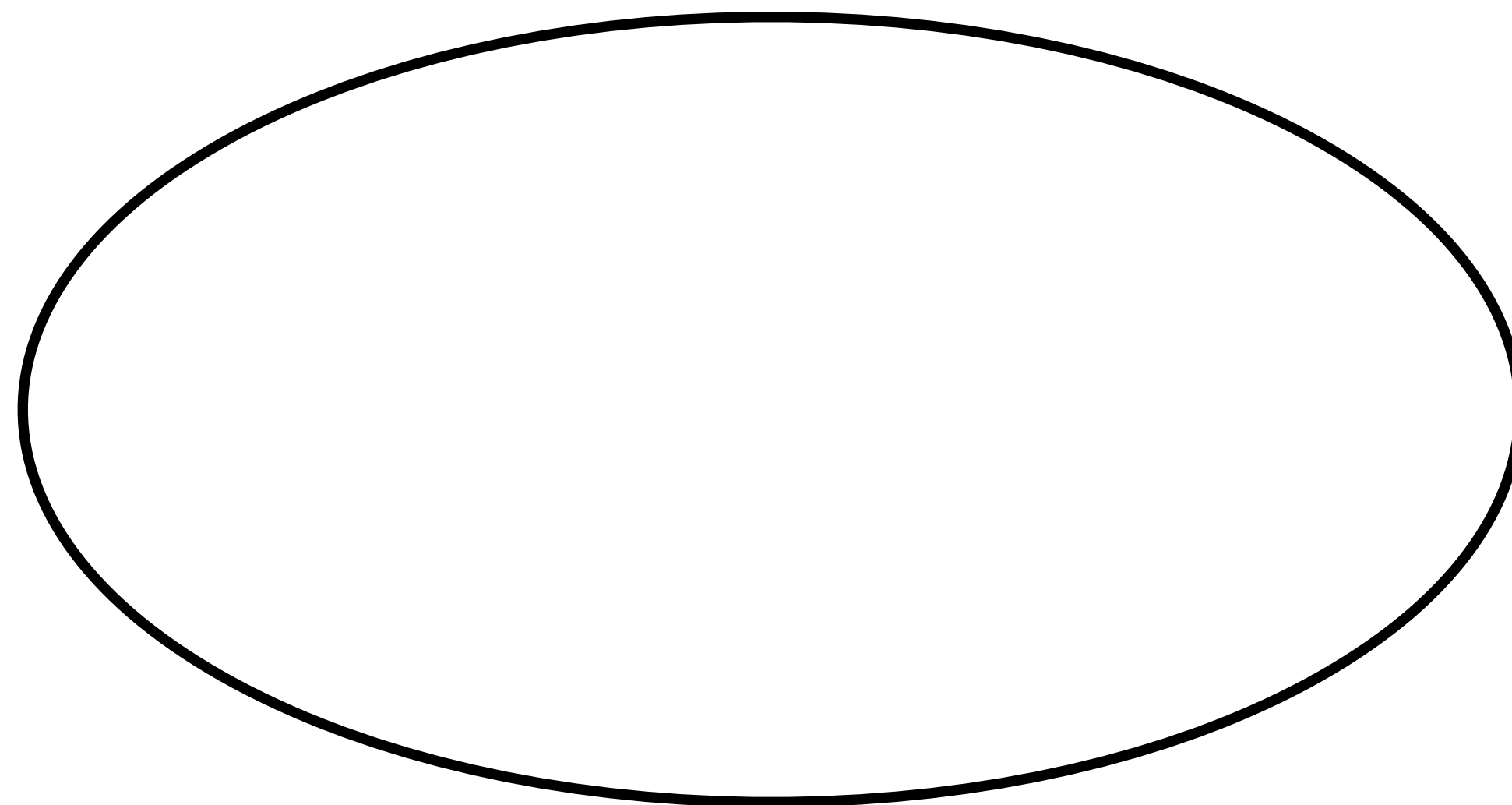
Vibe Coding

Review — Method

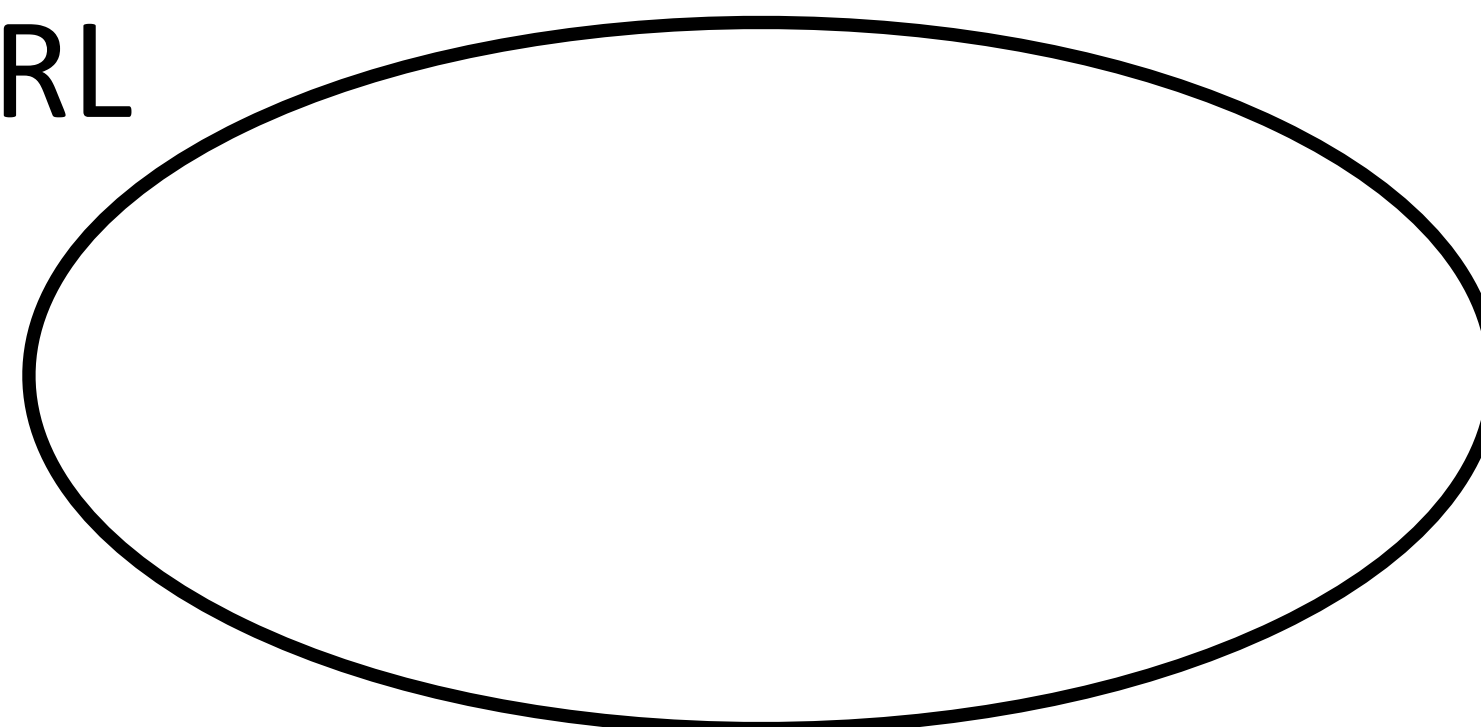
Pretraining



SFT

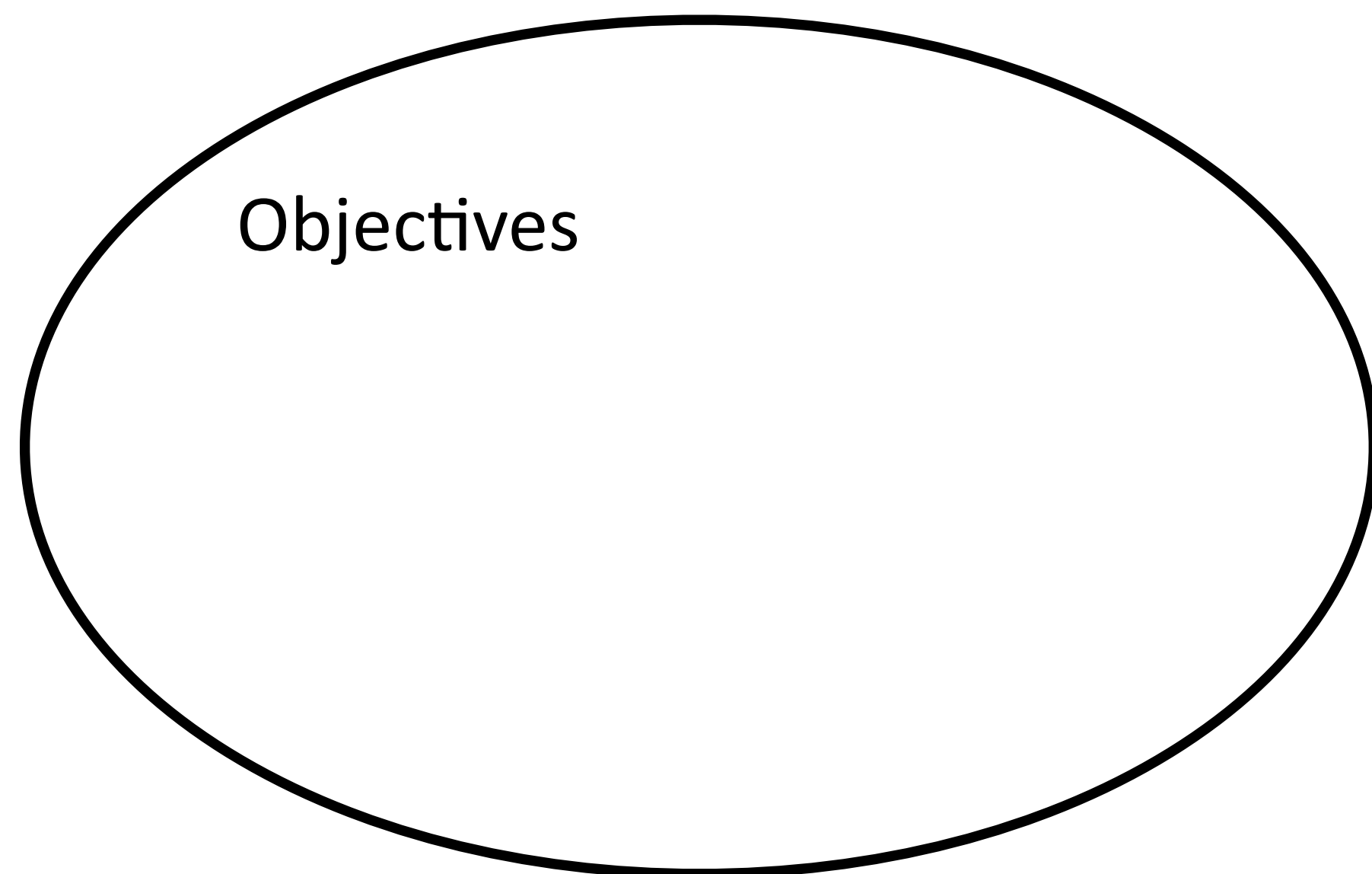


RL

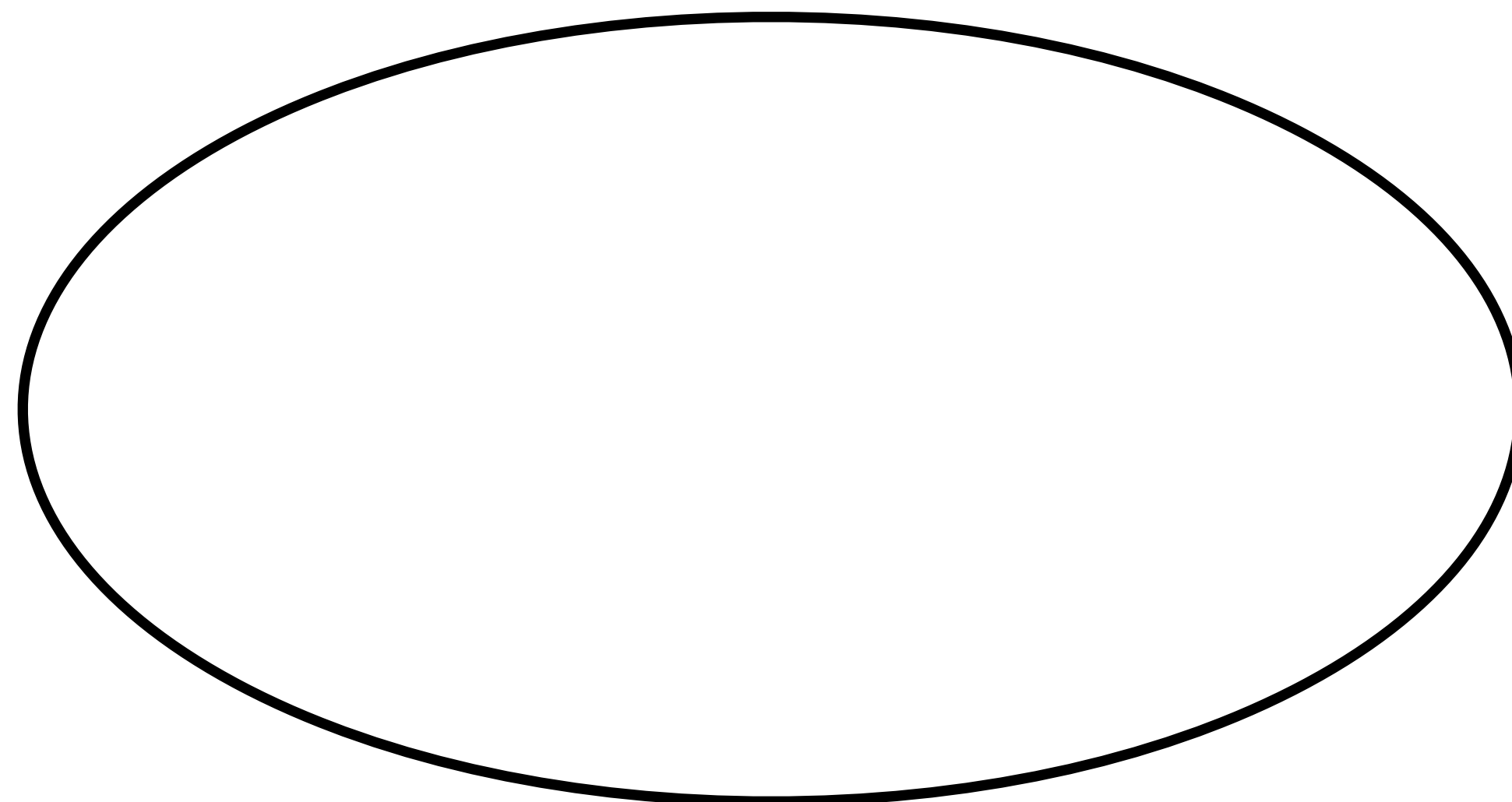


Review — Method

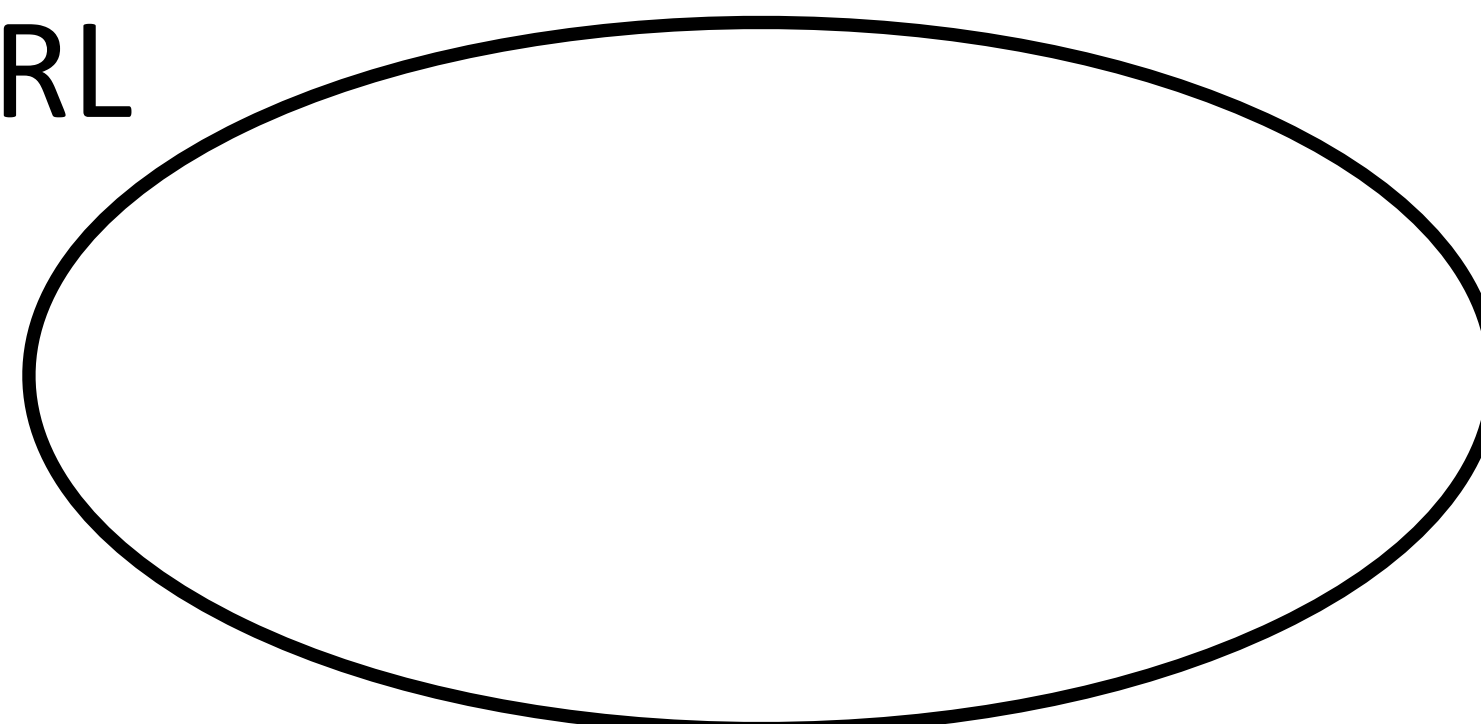
Pretraining



SFT

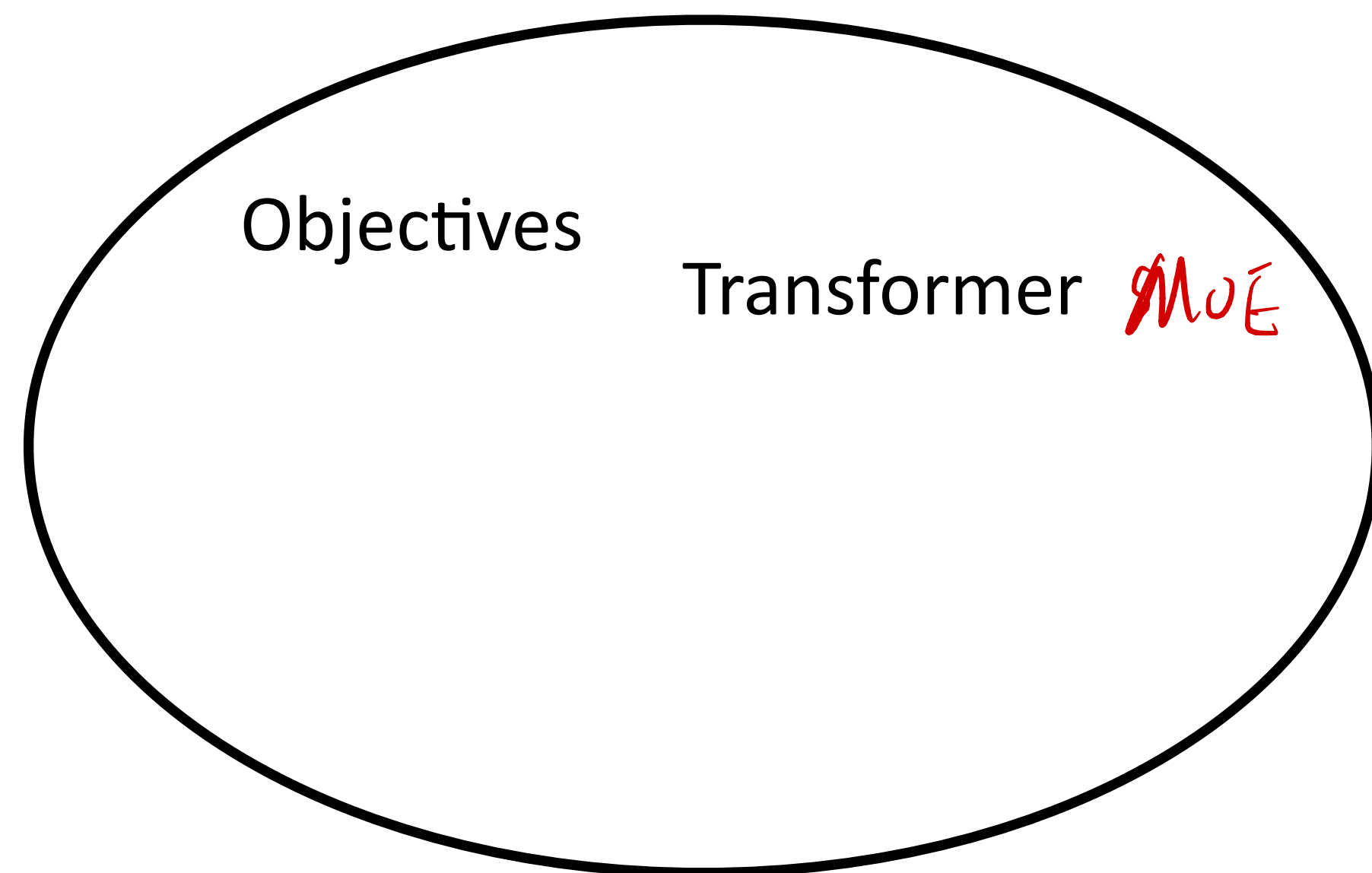


RL

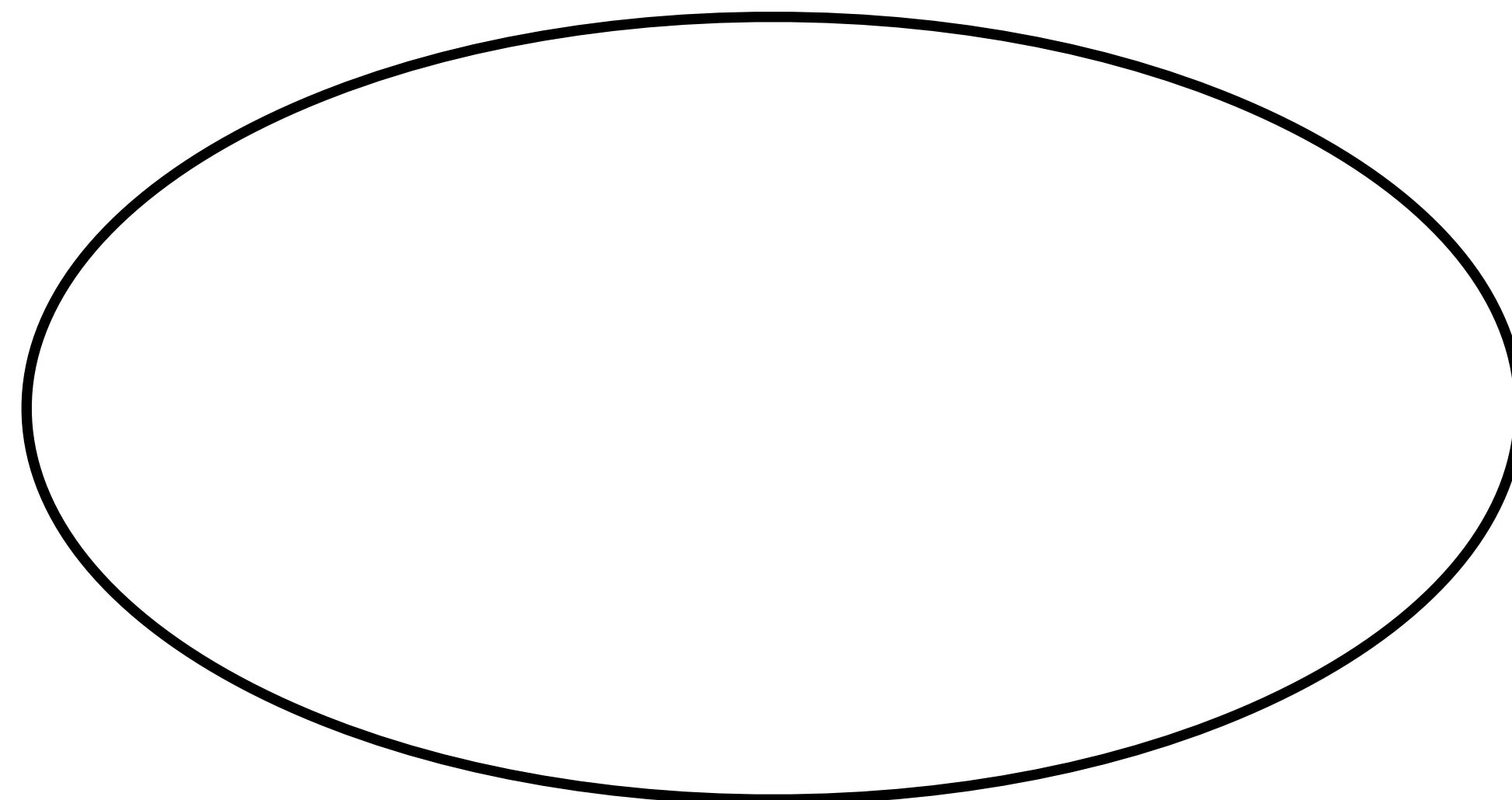


Review — Method

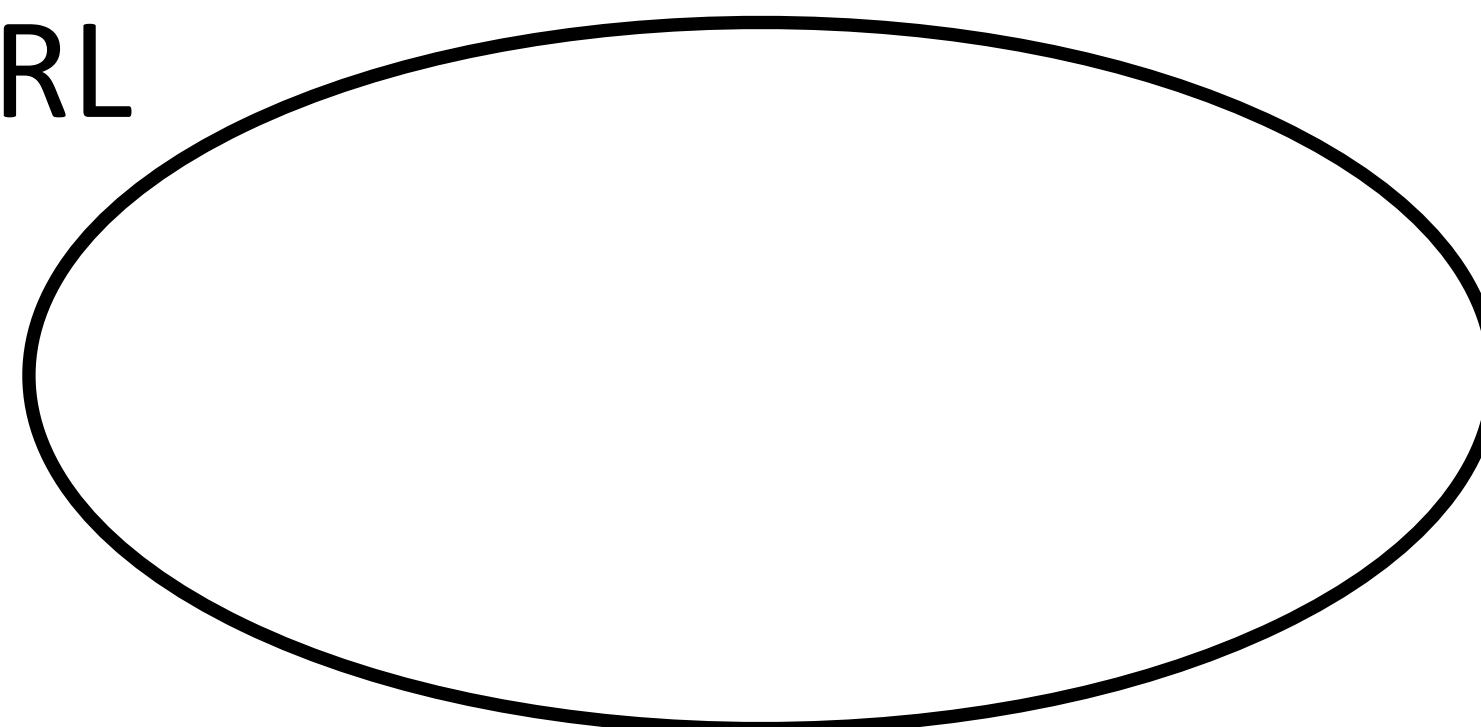
Pretraining



SFT

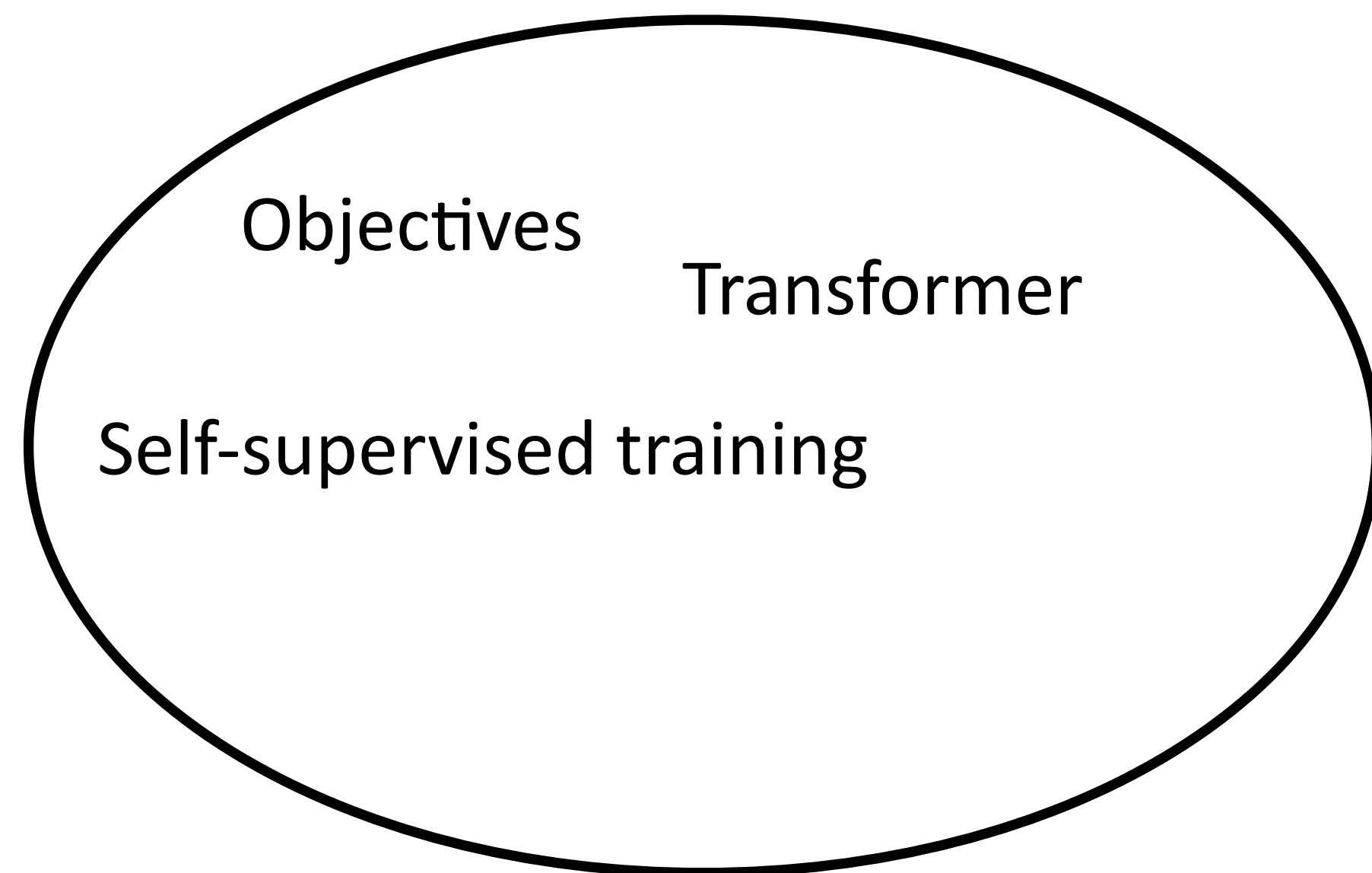


RL

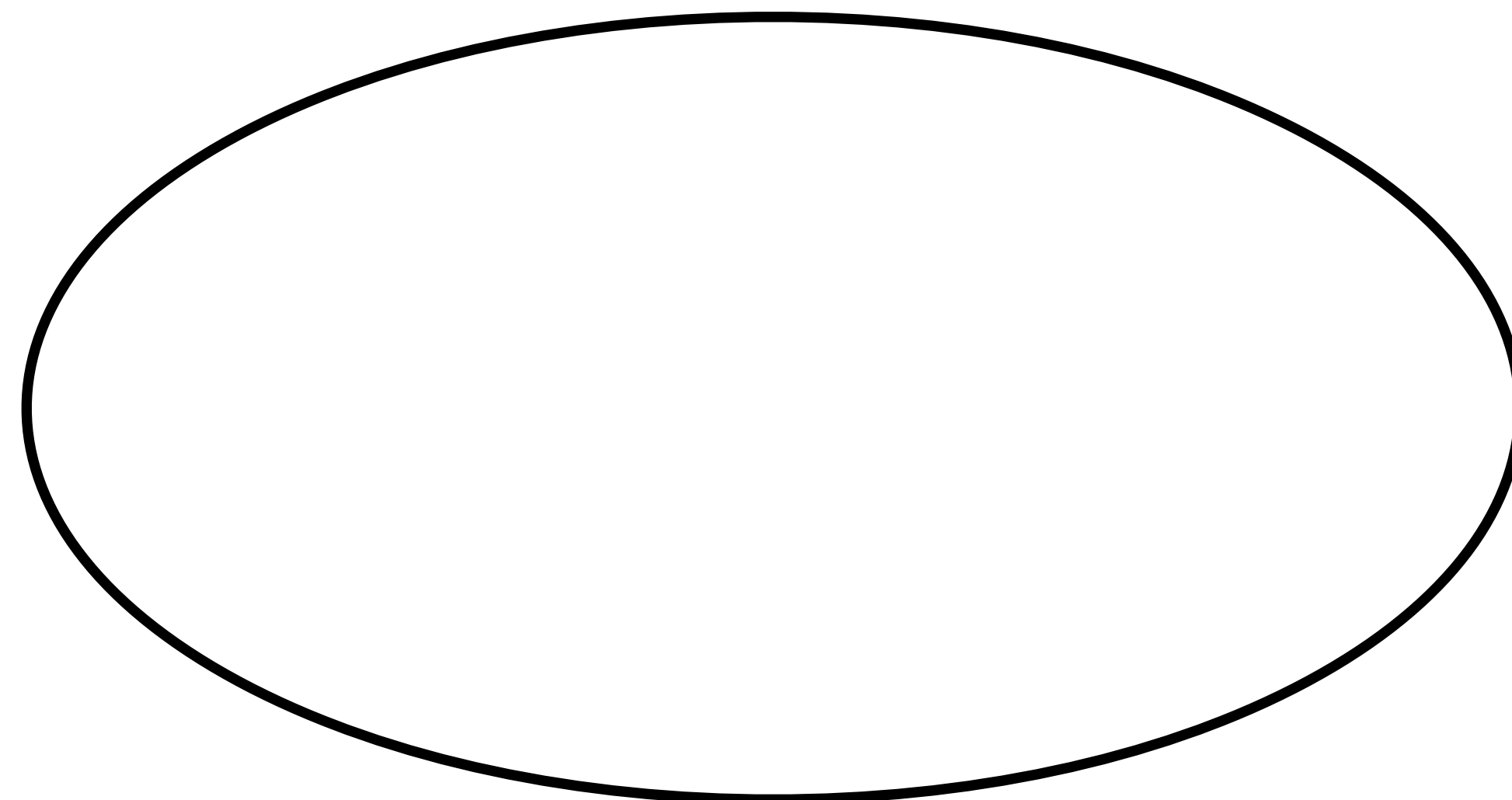


Review — Method

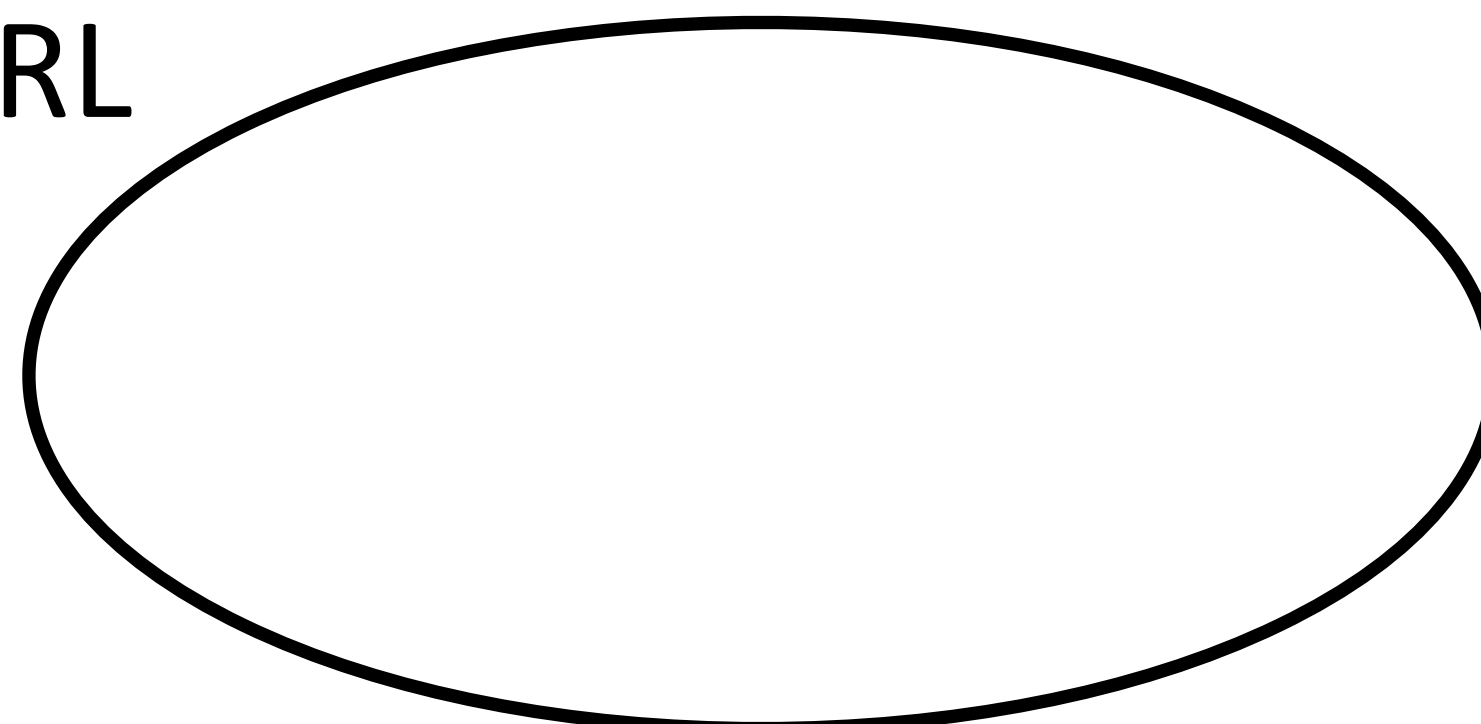
Pretraining



SFT

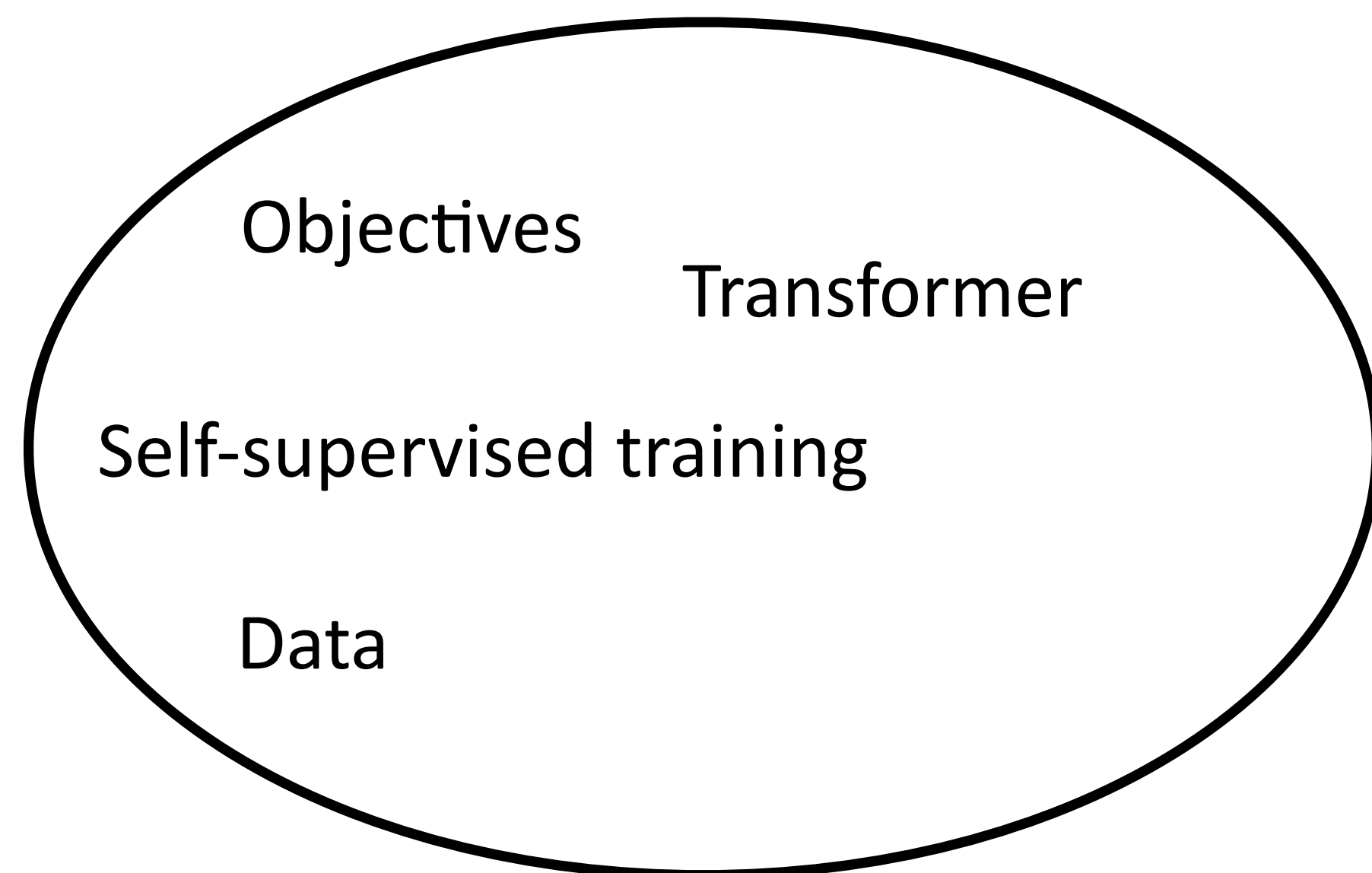


RL

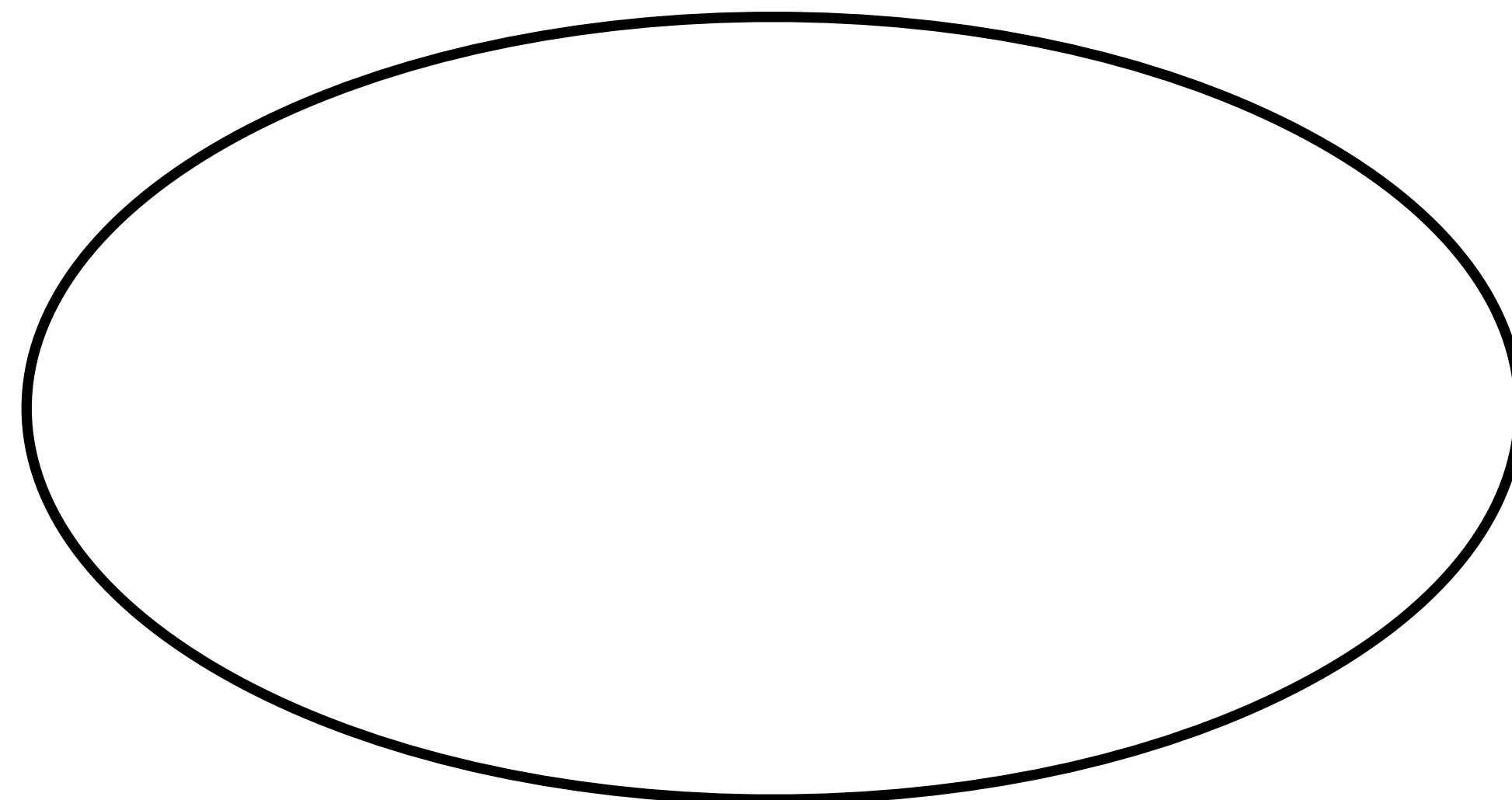


Review — Method

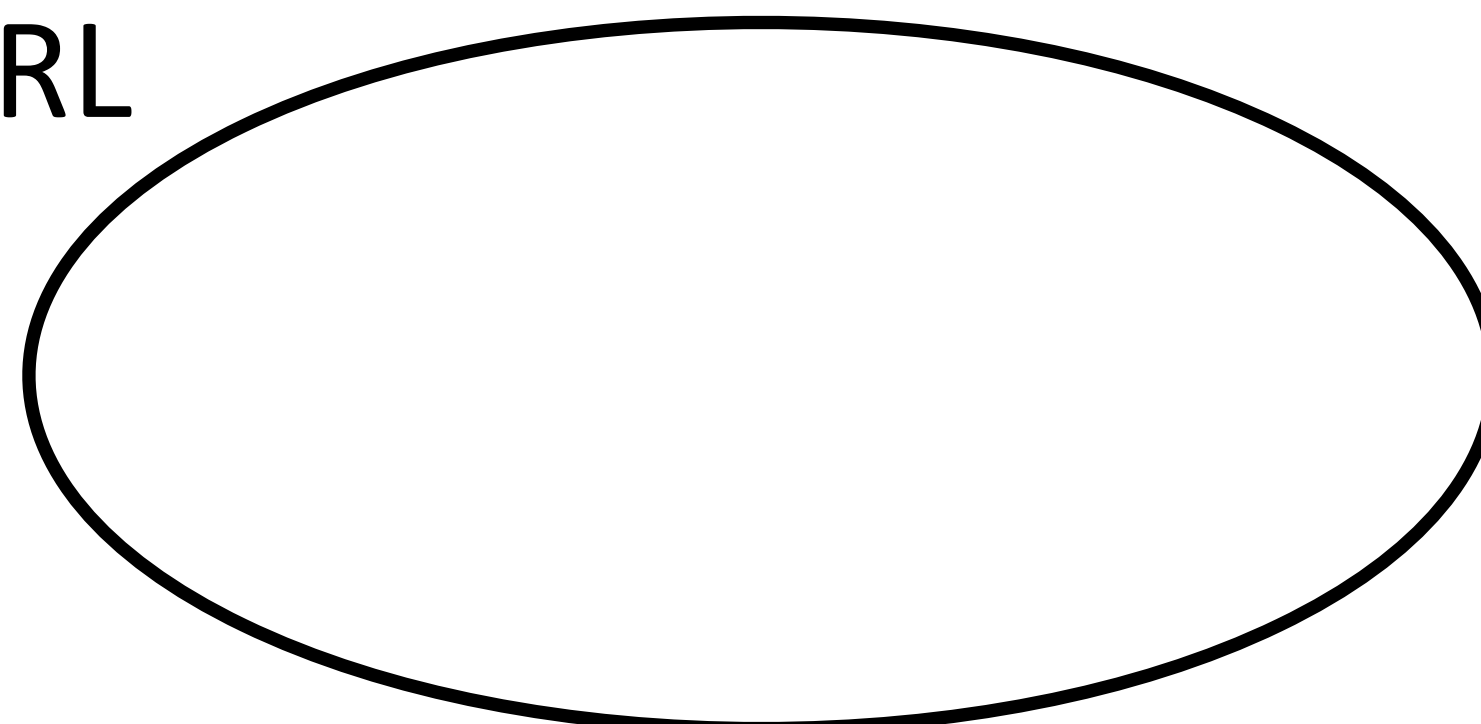
Pretraining



SFT

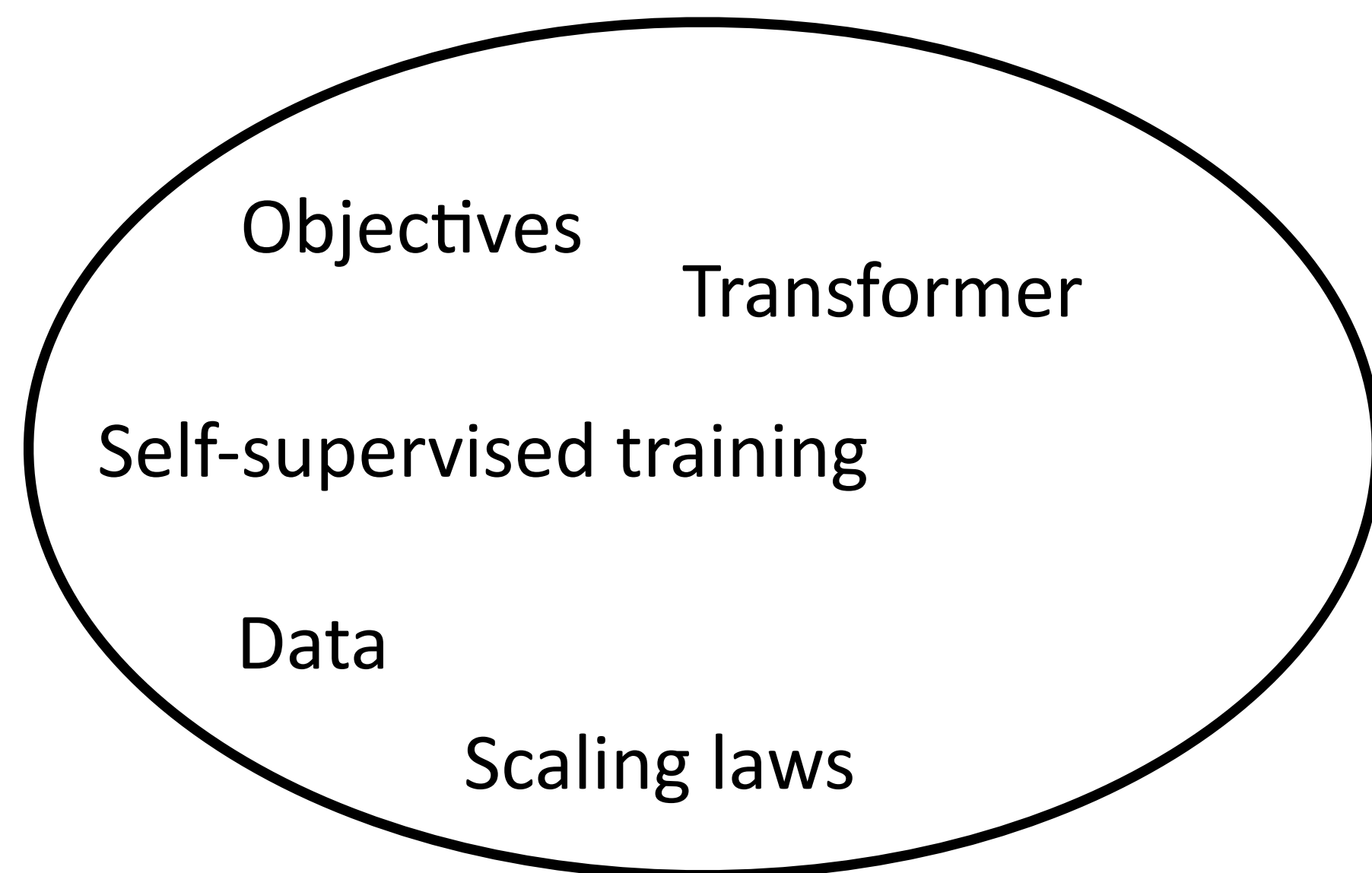


RL

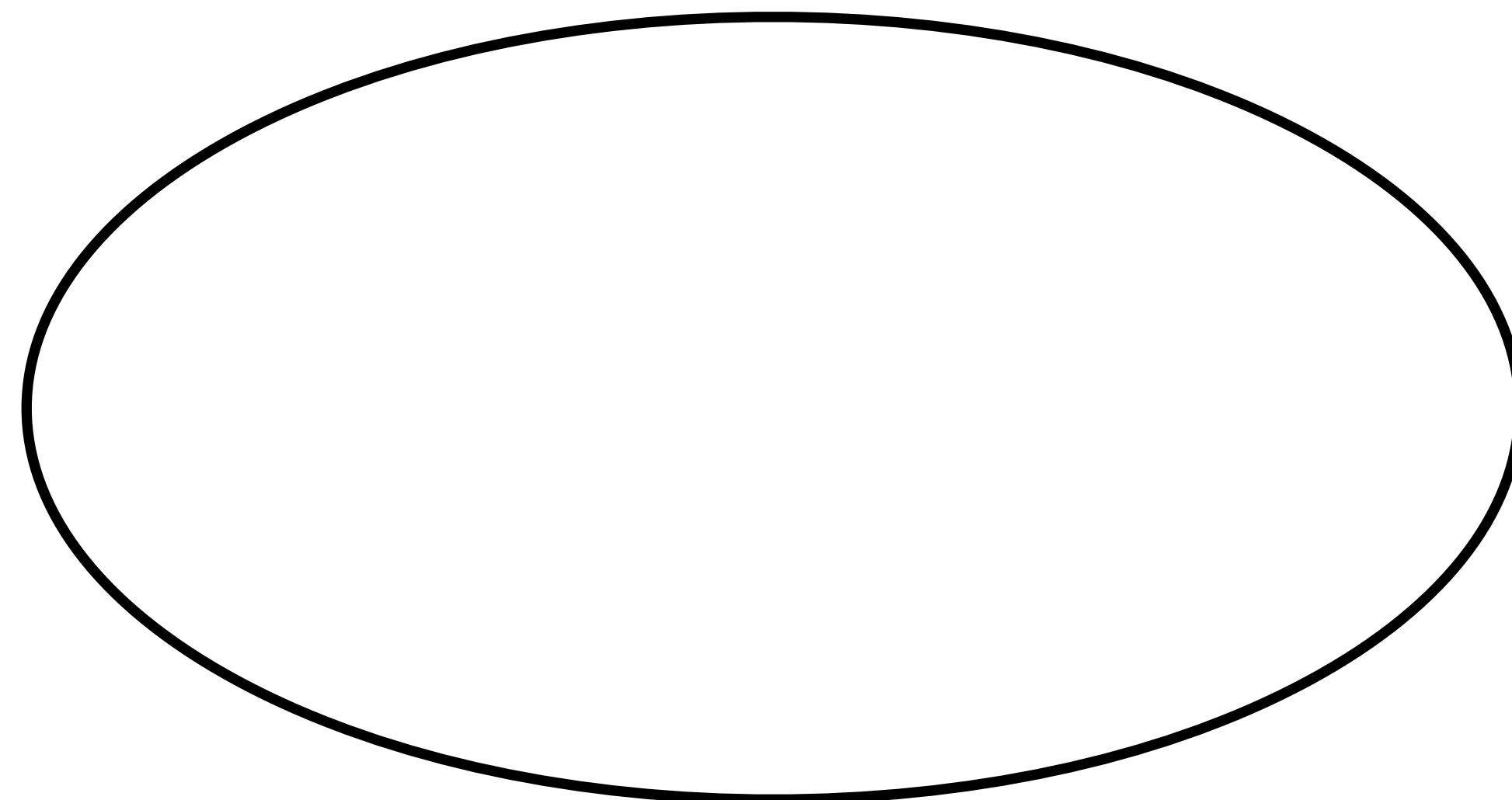


Review — Method

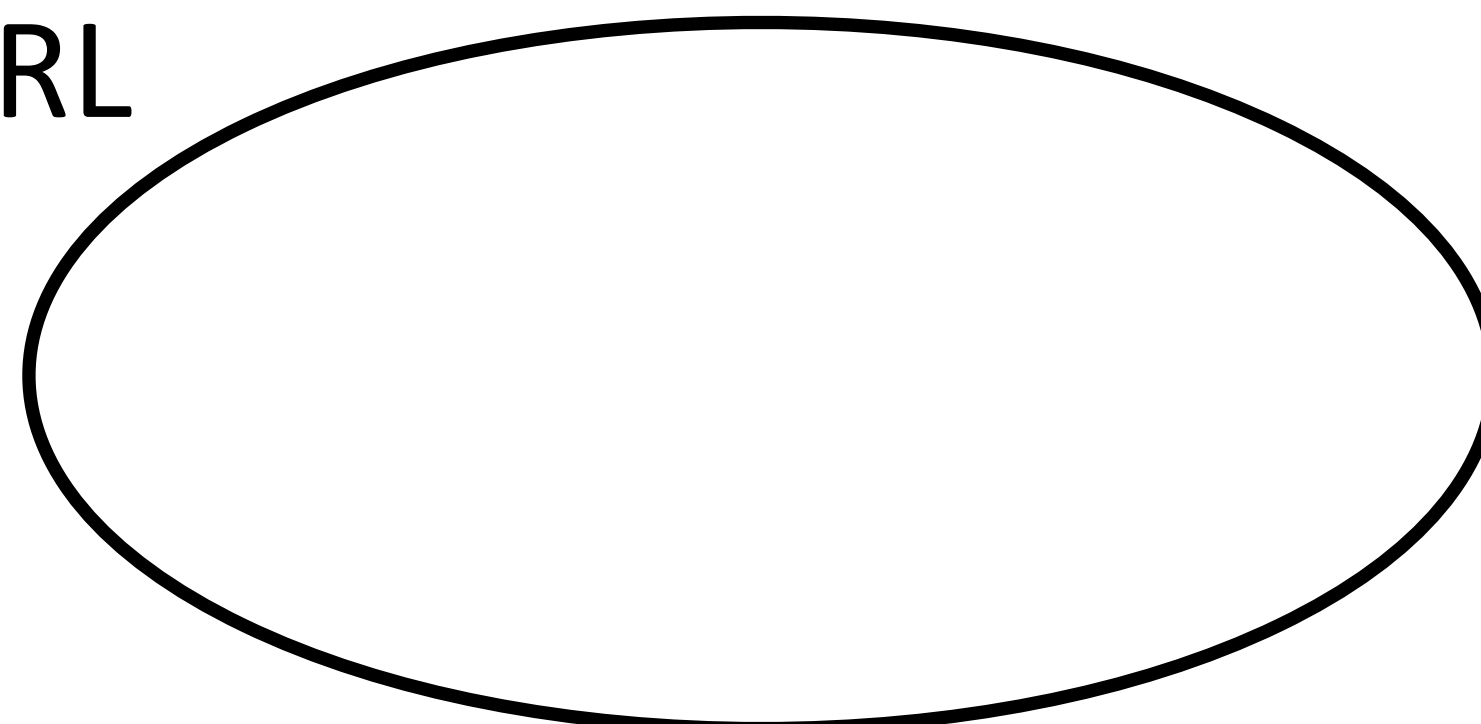
Pretraining



SFT

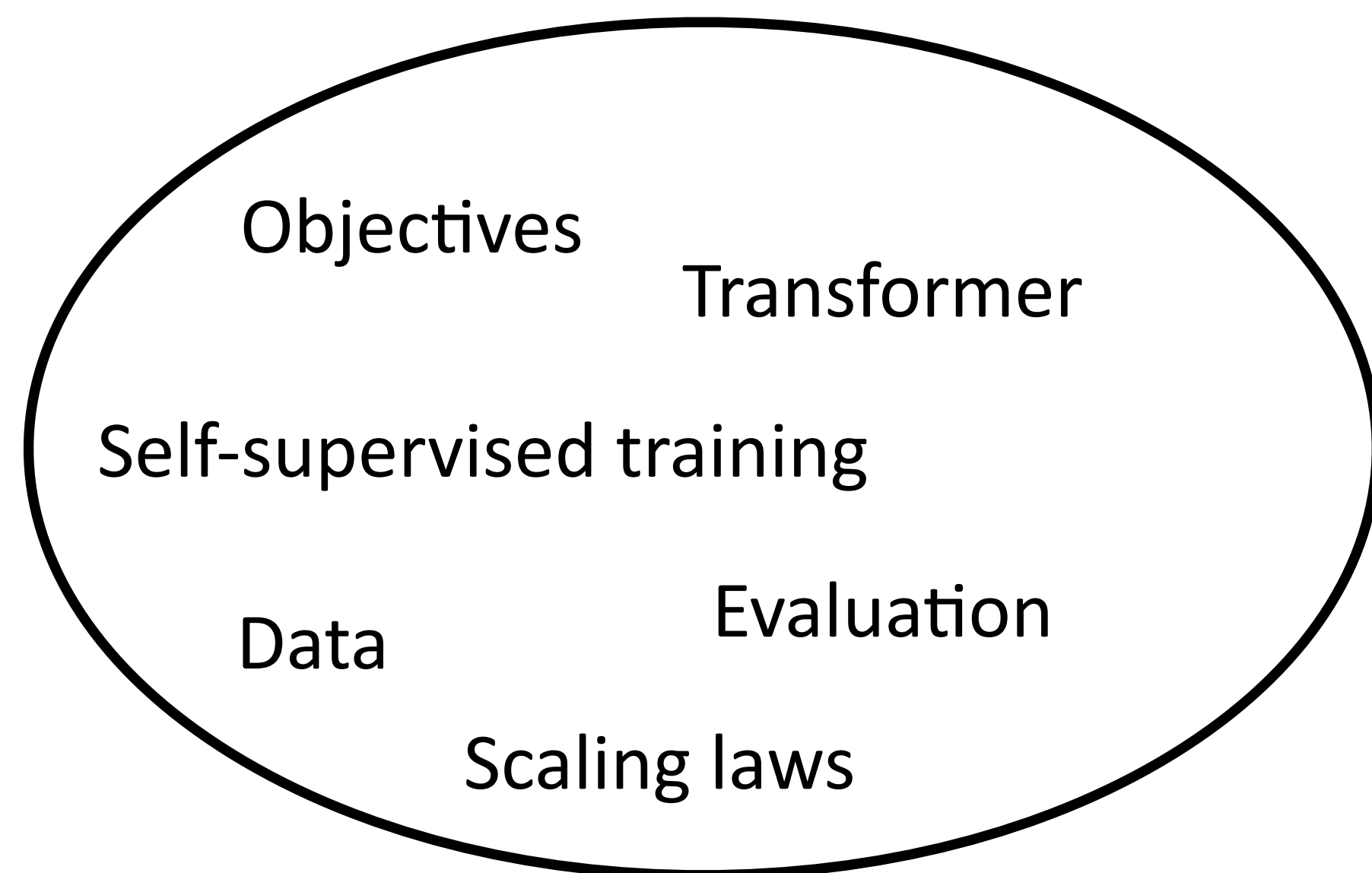


RL

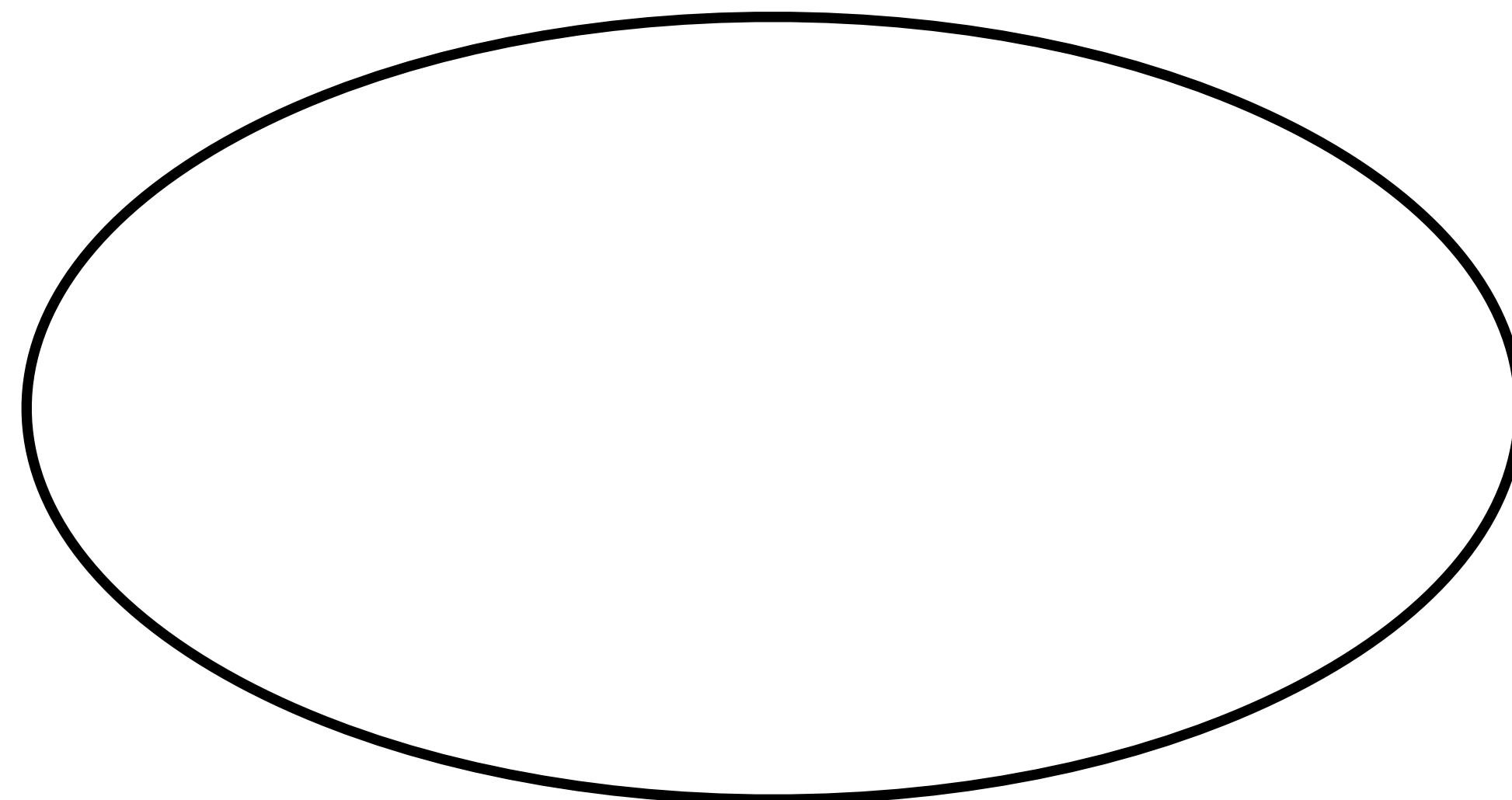


Review — Method

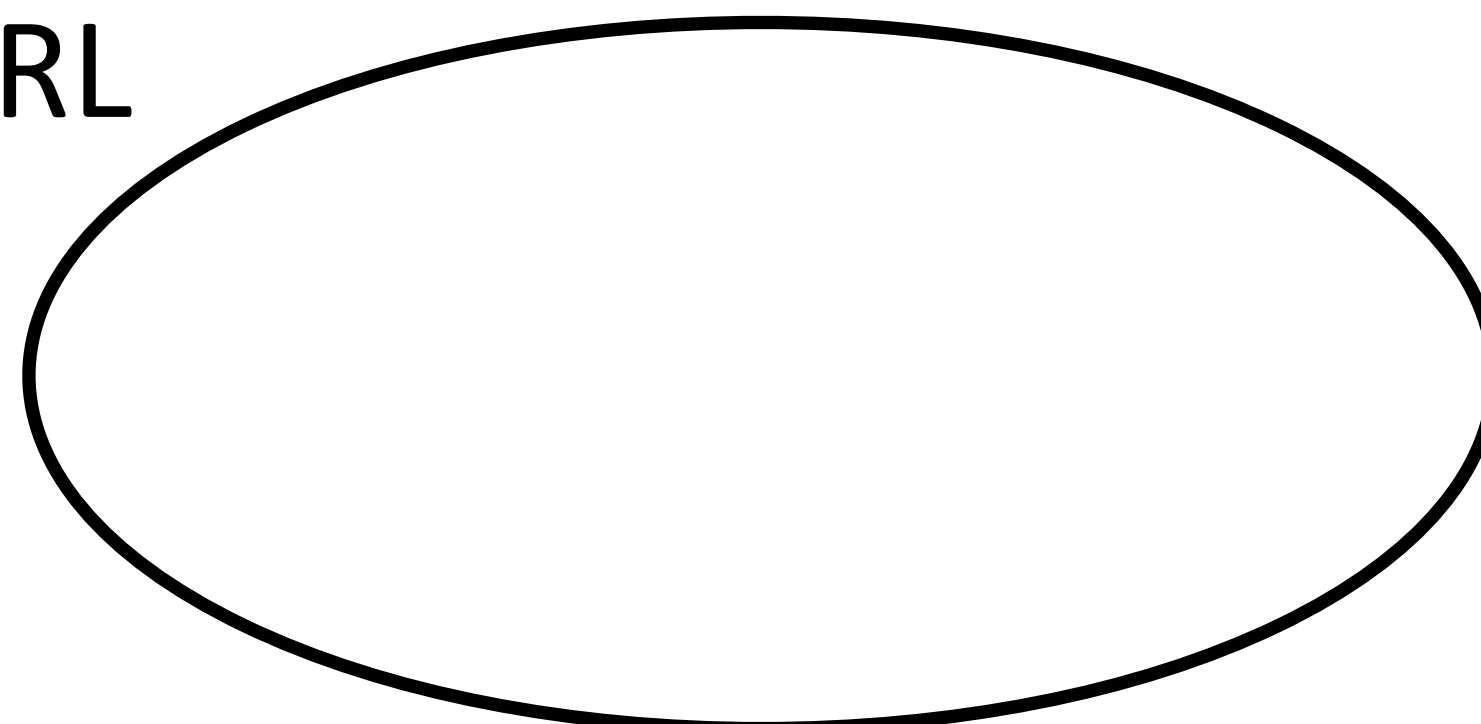
Pretraining



SFT

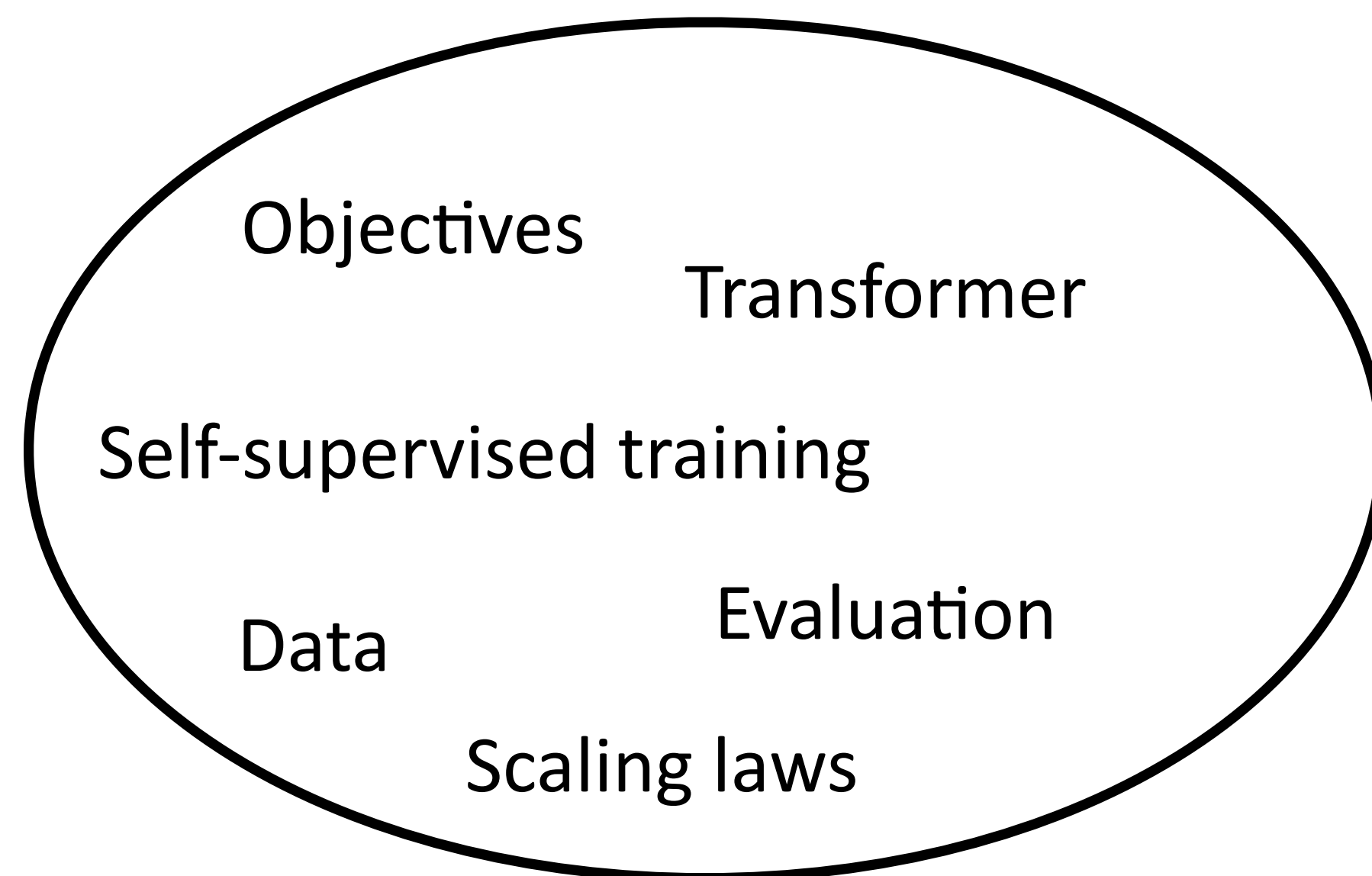


RL

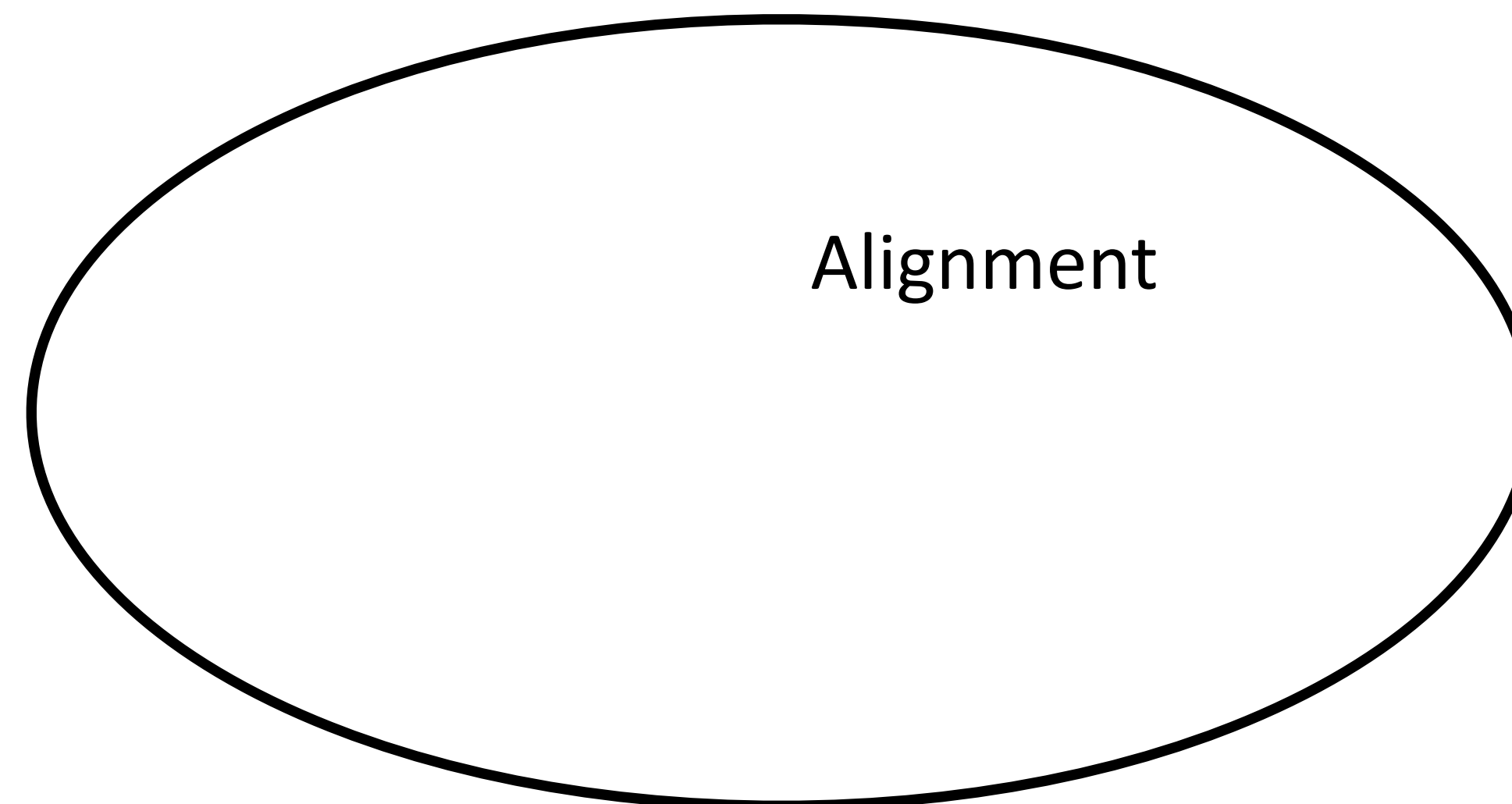


Review — Method

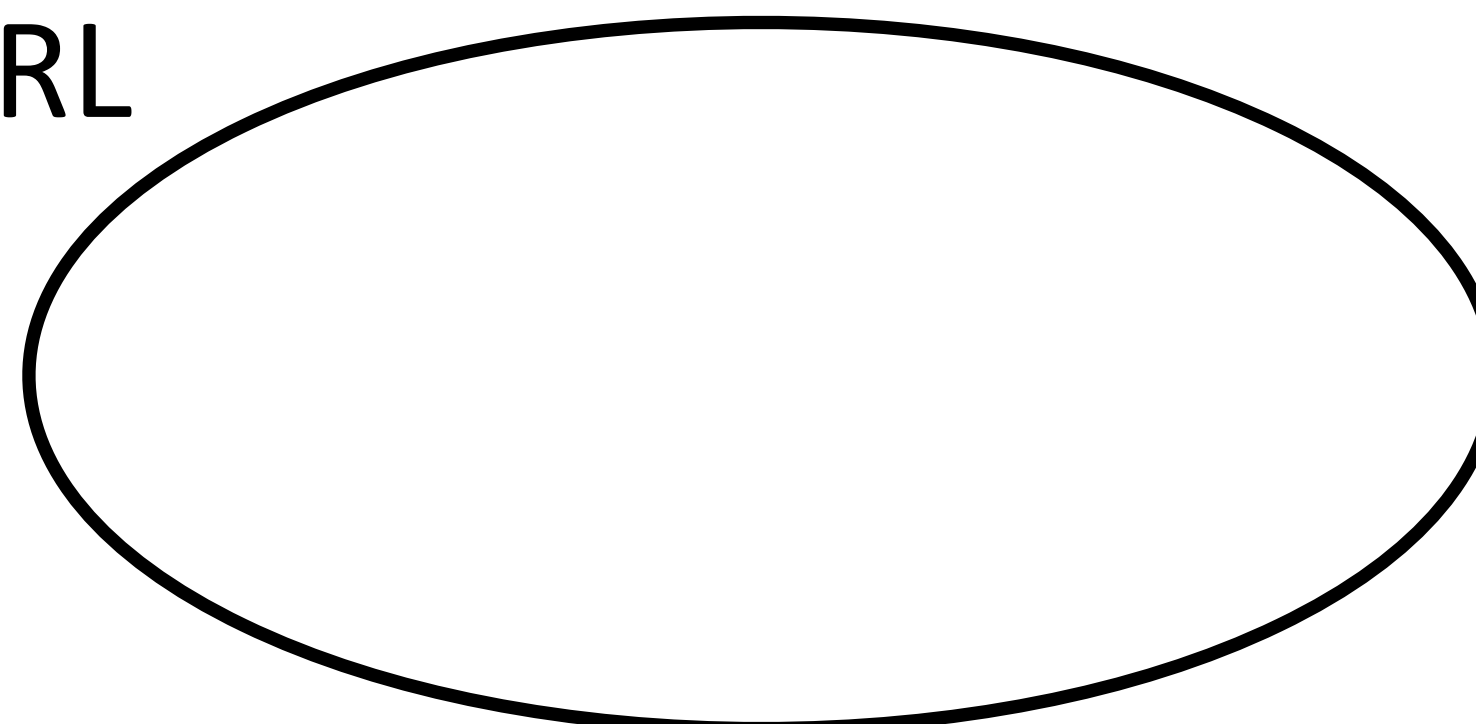
Pretraining



SFT

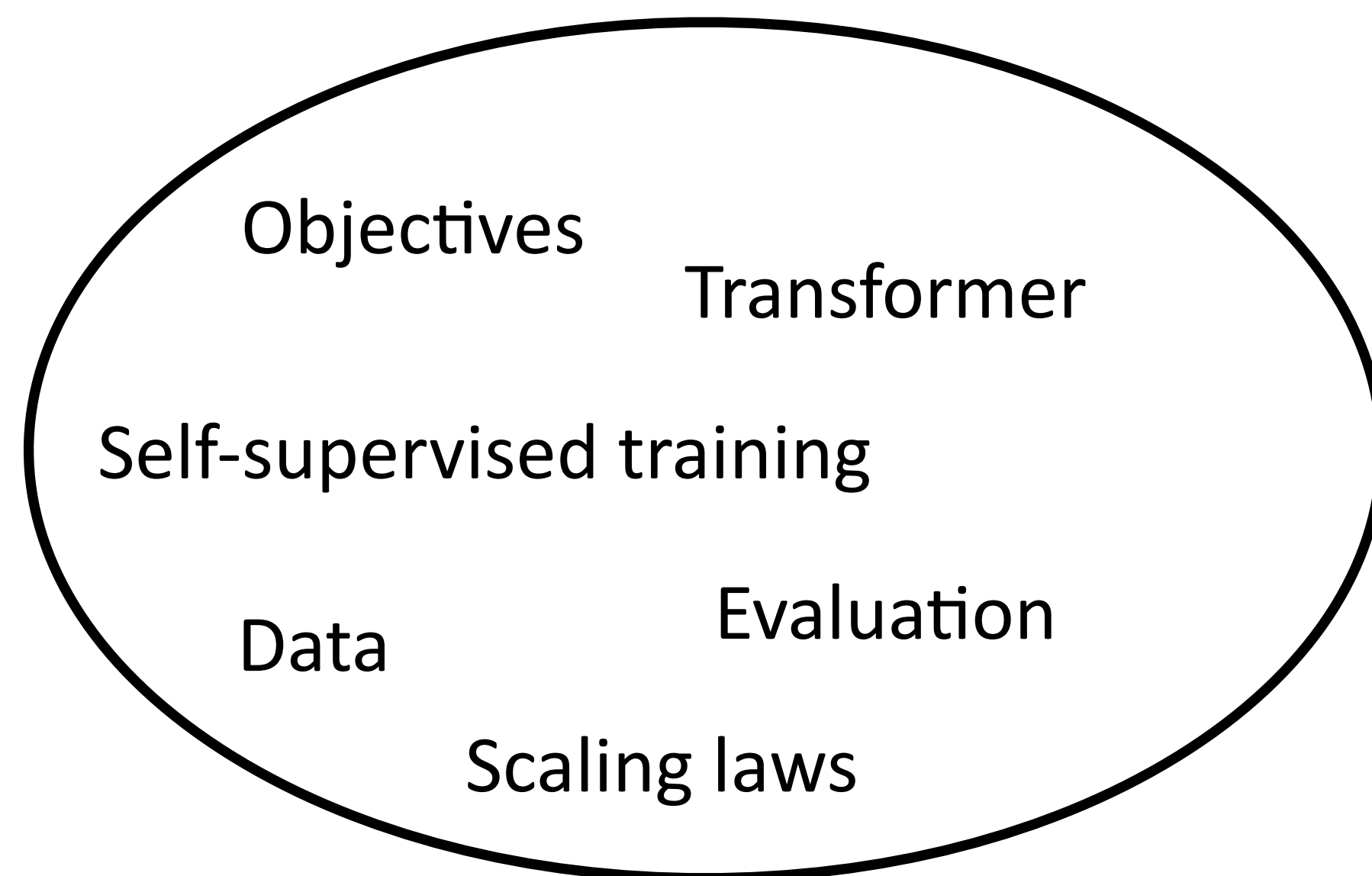


RL

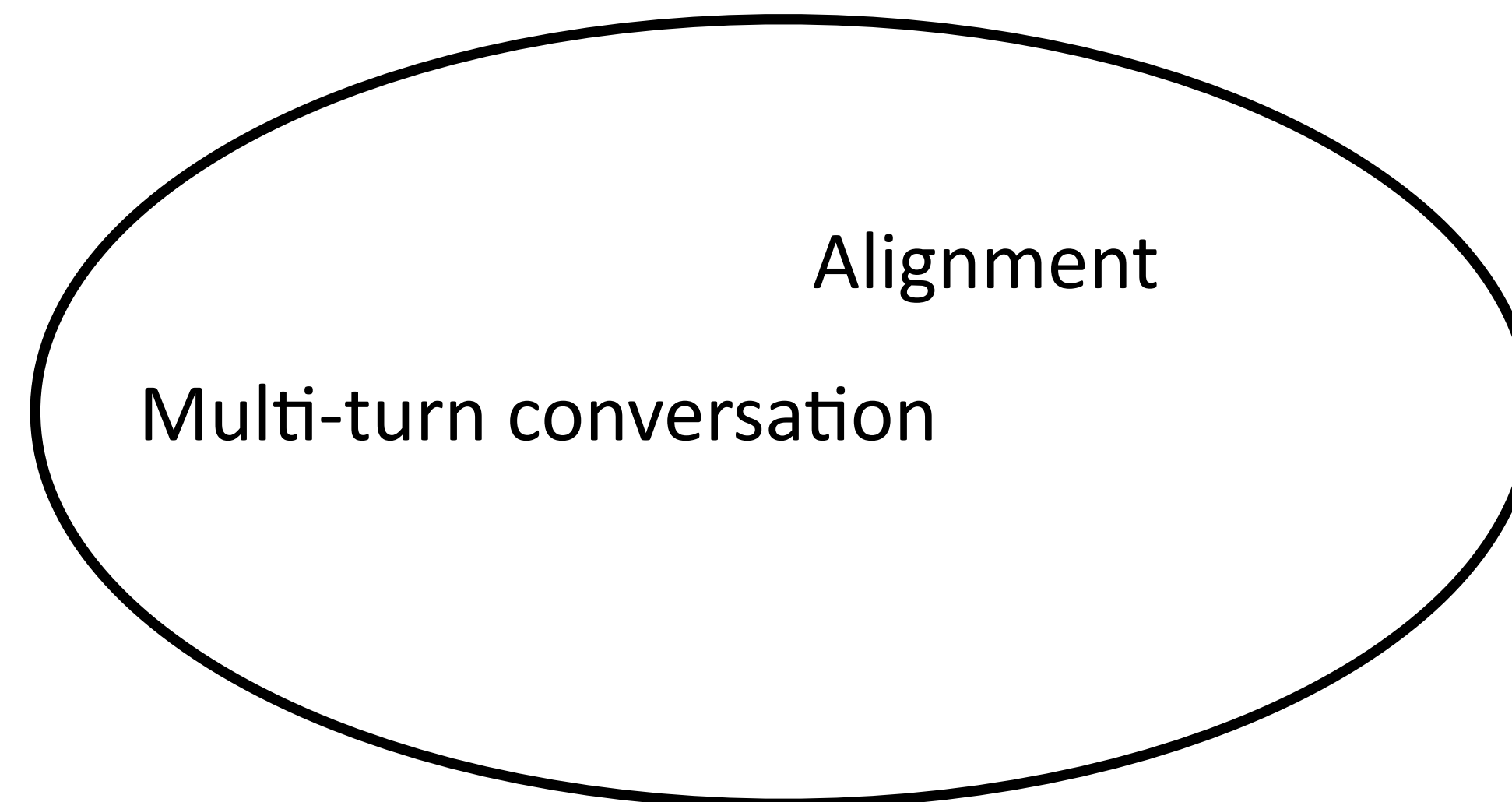


Review — Method

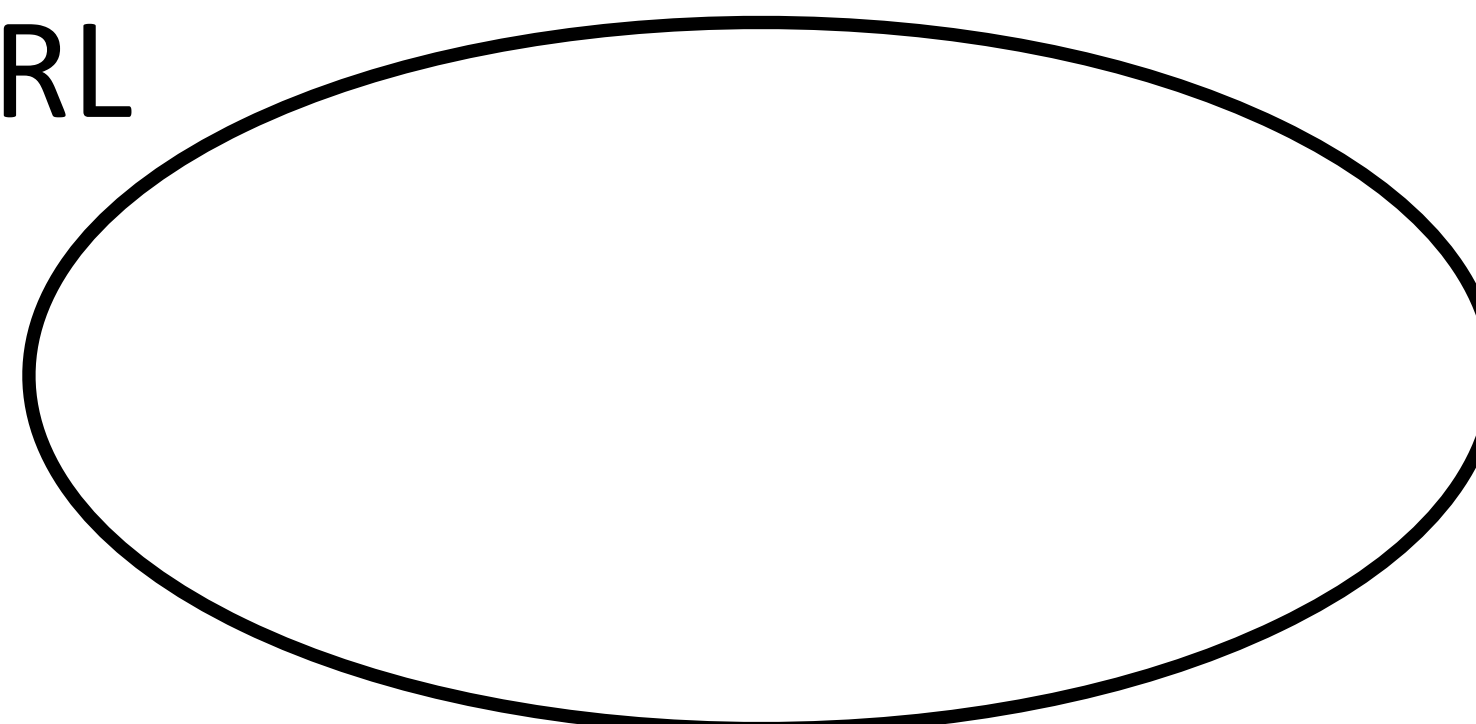
Pretraining



SFT

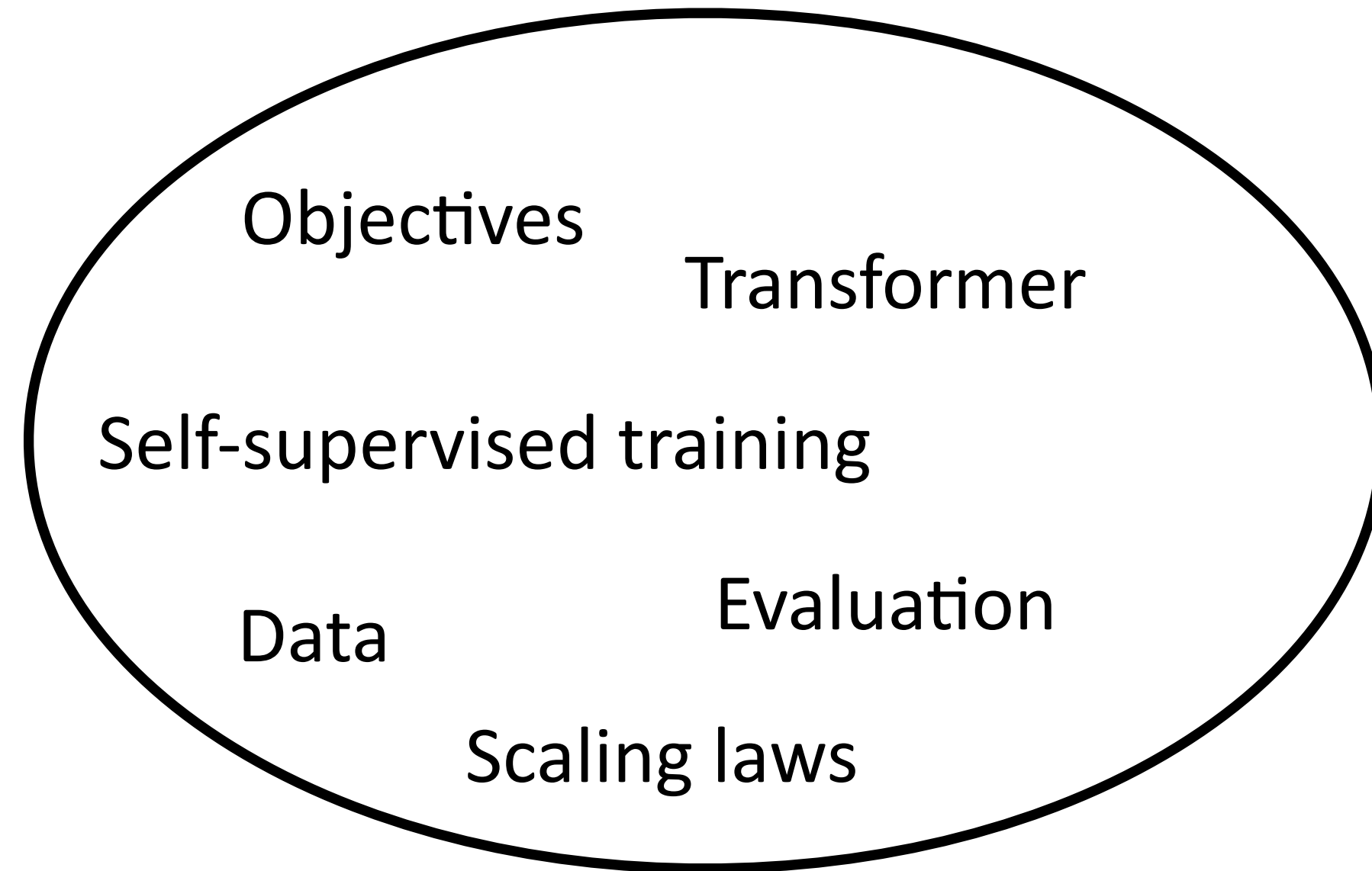


RL

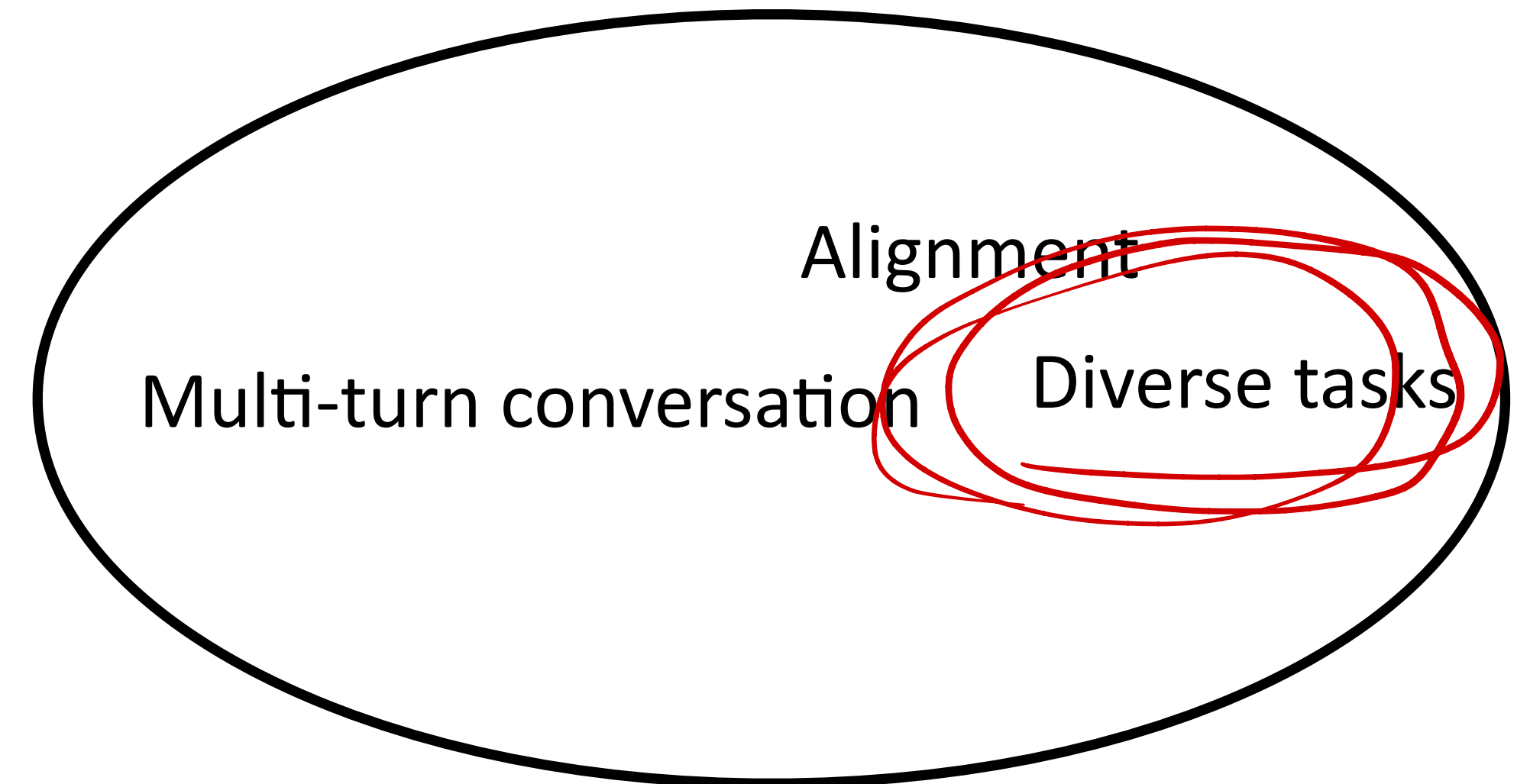


Review — Method

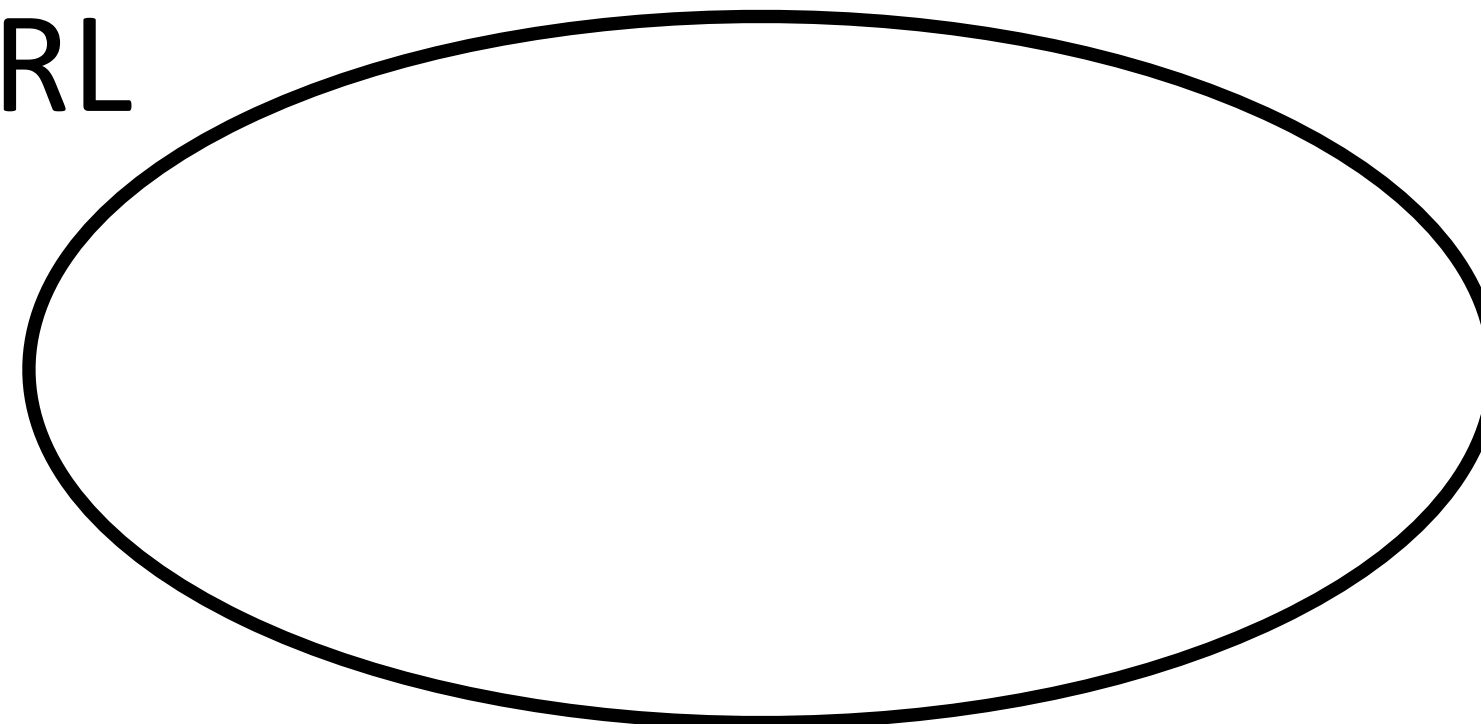
Pretraining



SFT

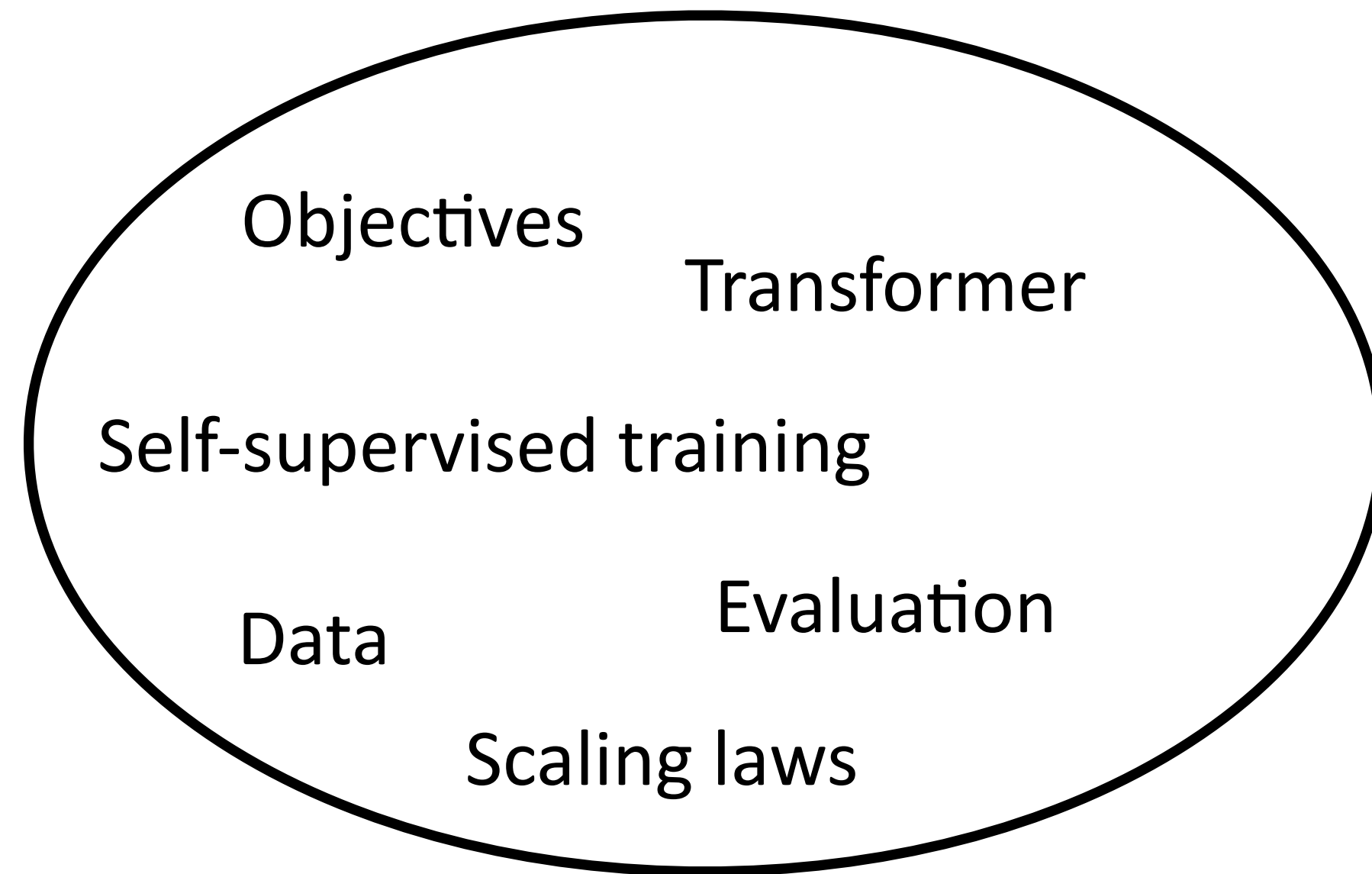


RL



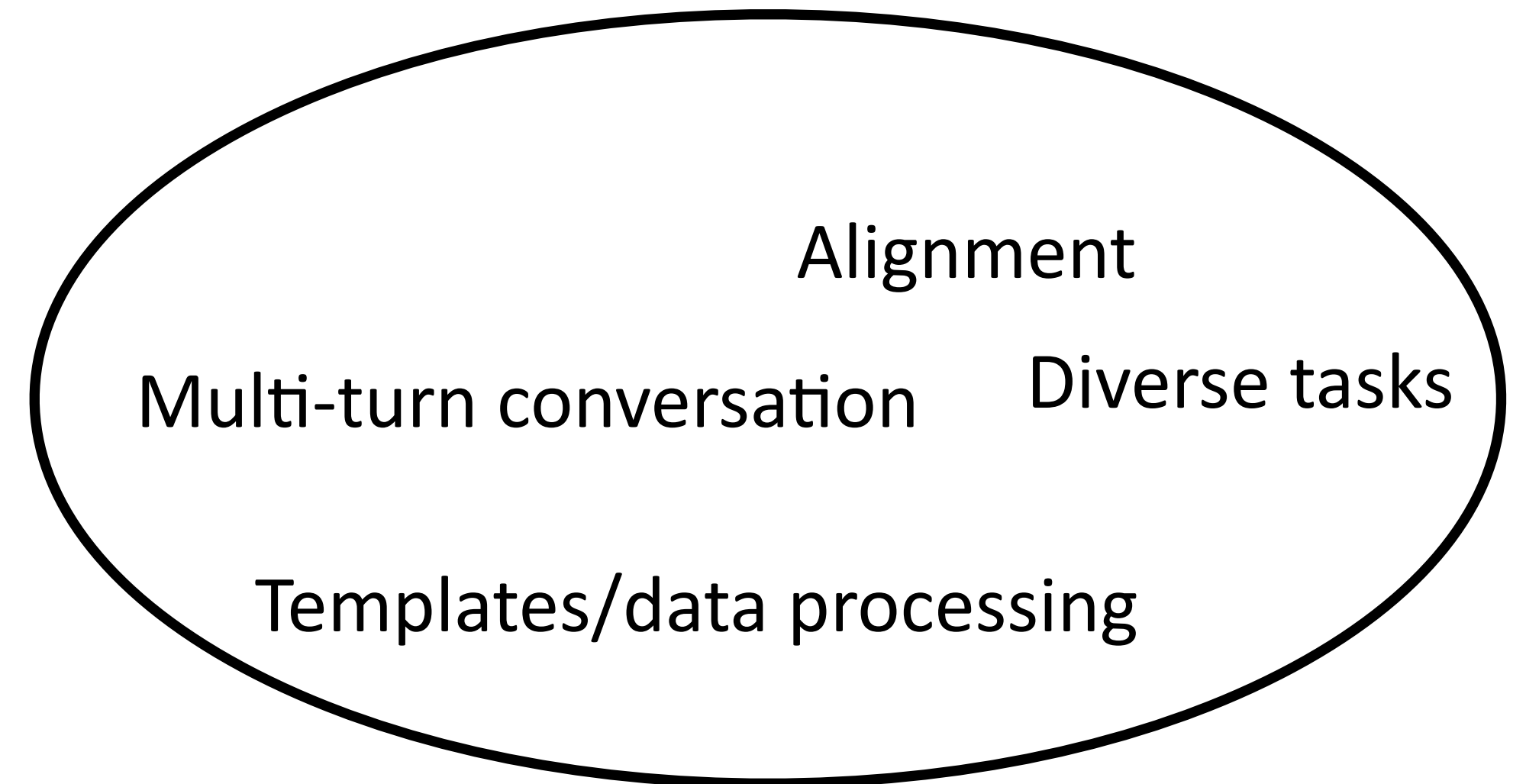
Review — Method

Pretraining

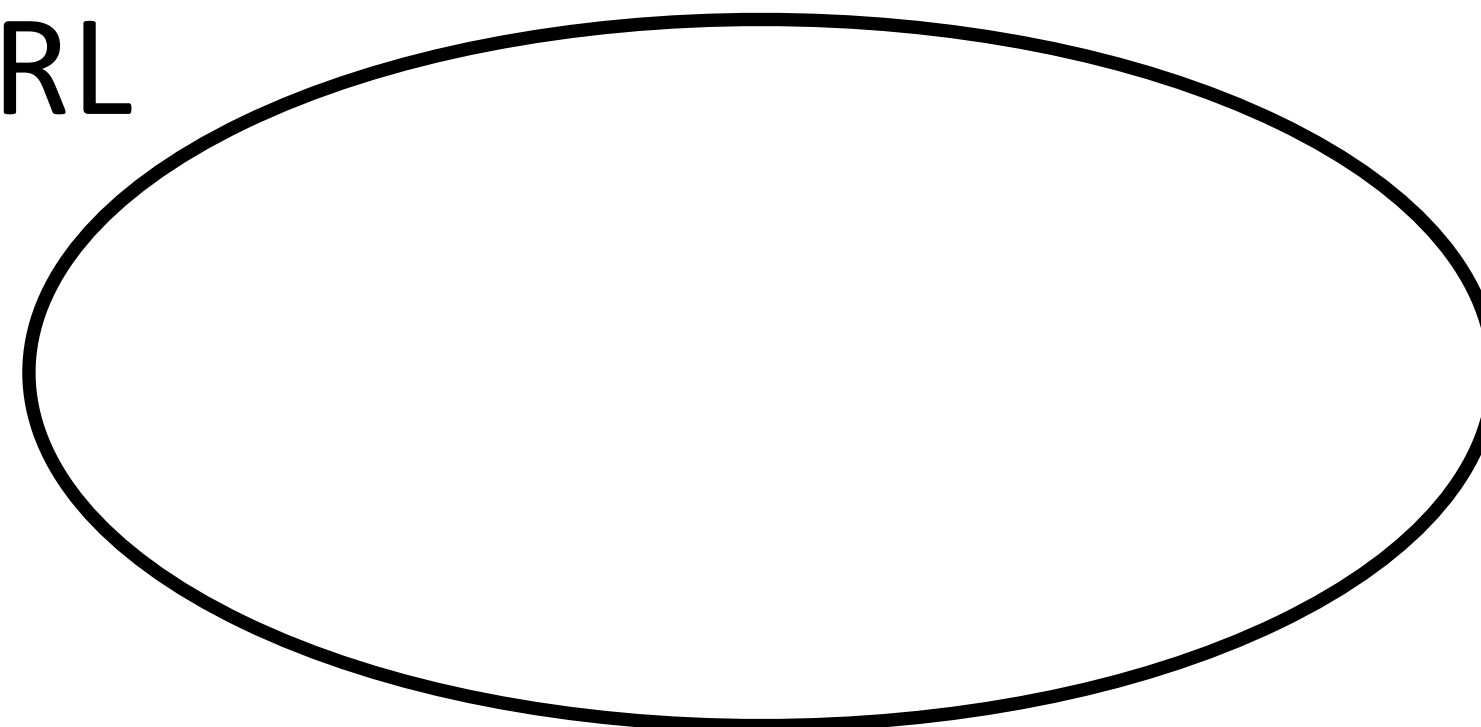


SFT

= instruction tuning

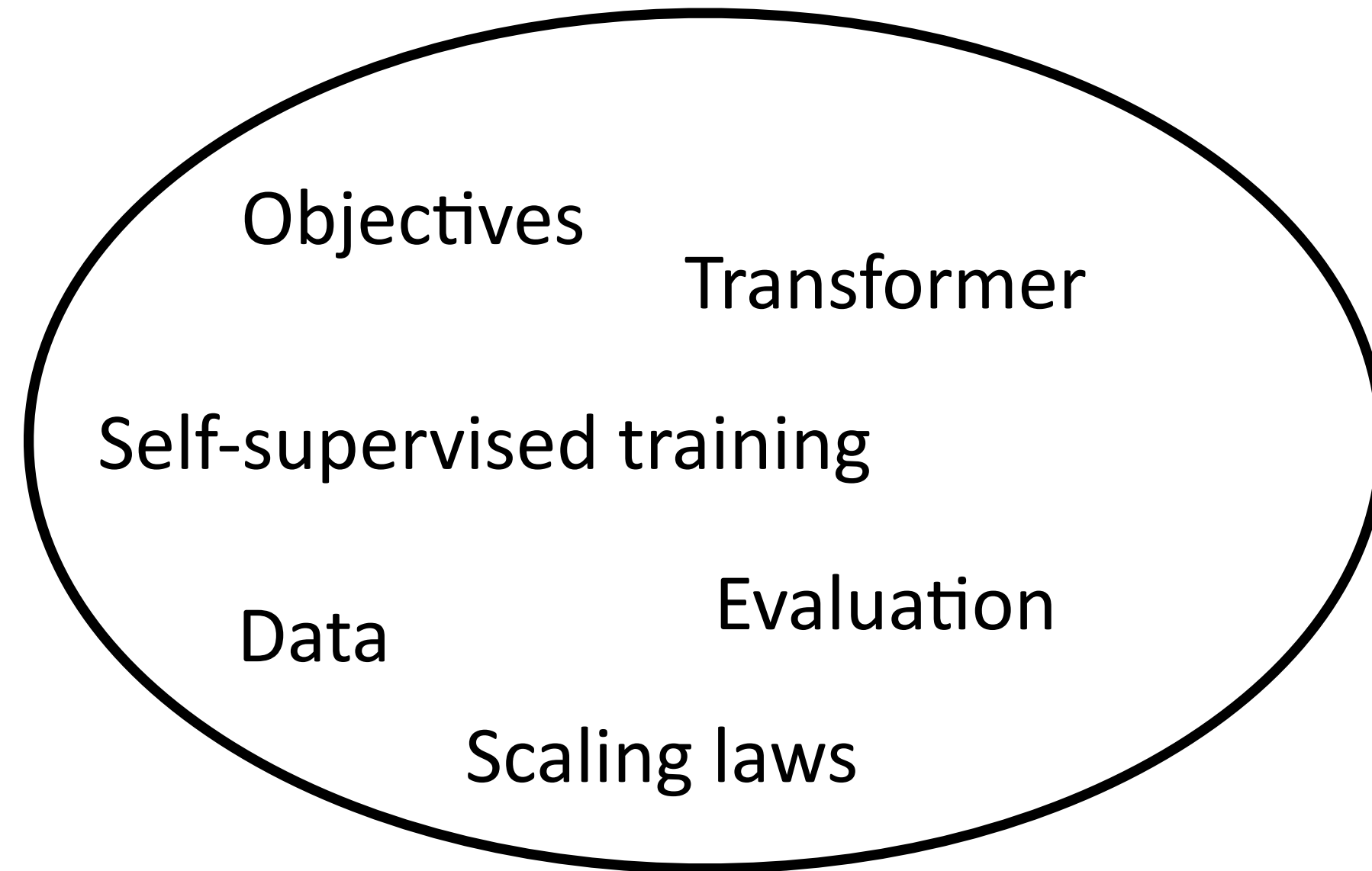


RL

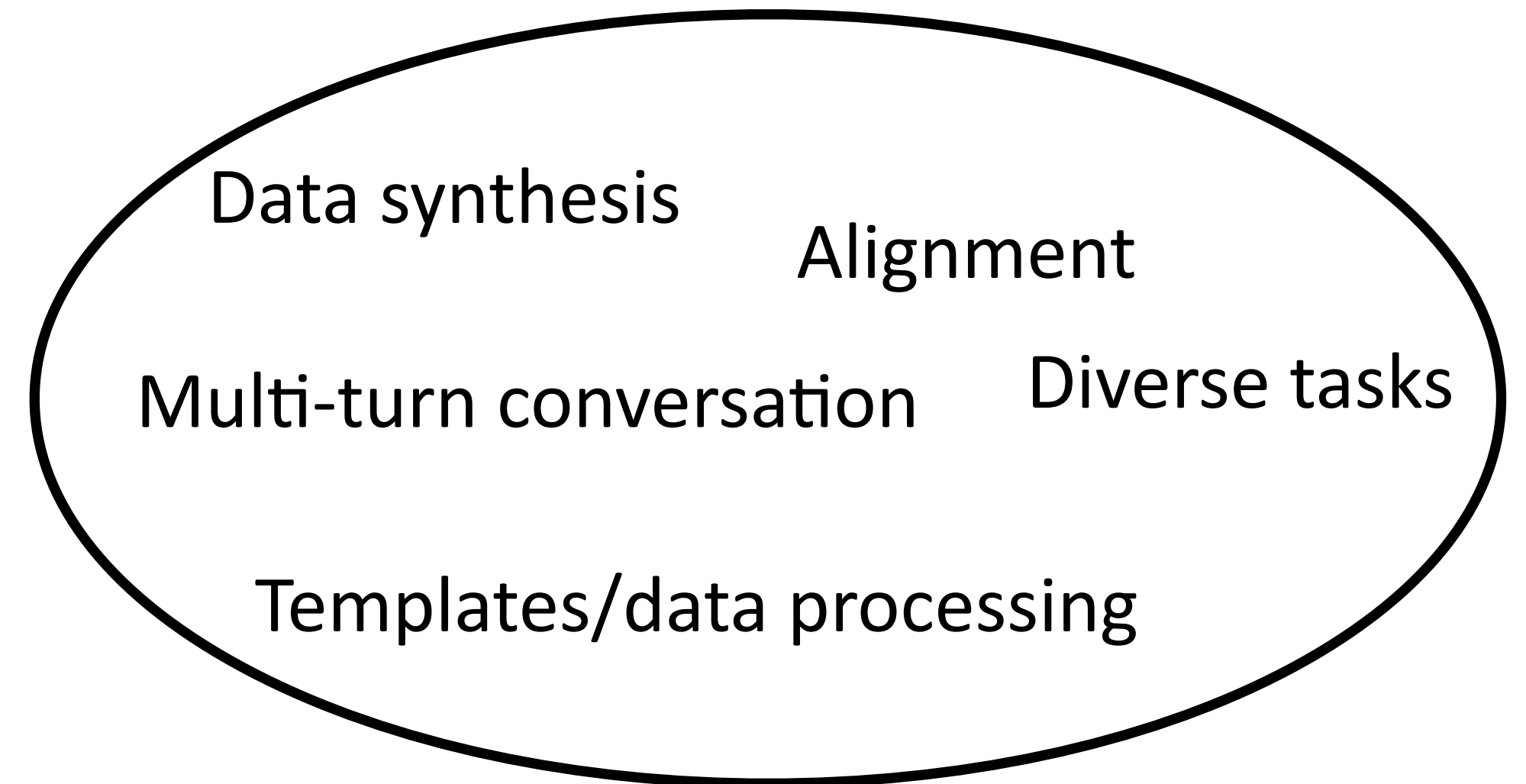


Review — Method

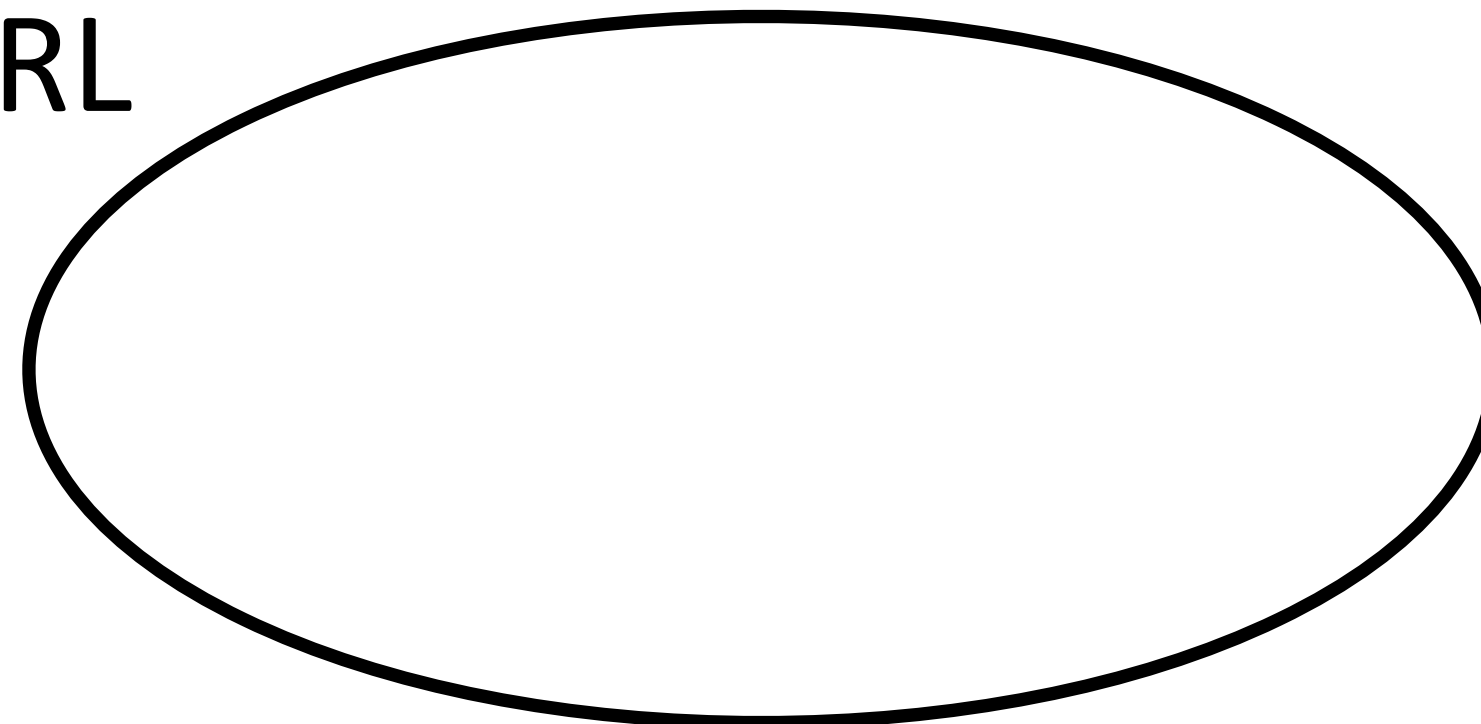
Pretraining



SFT

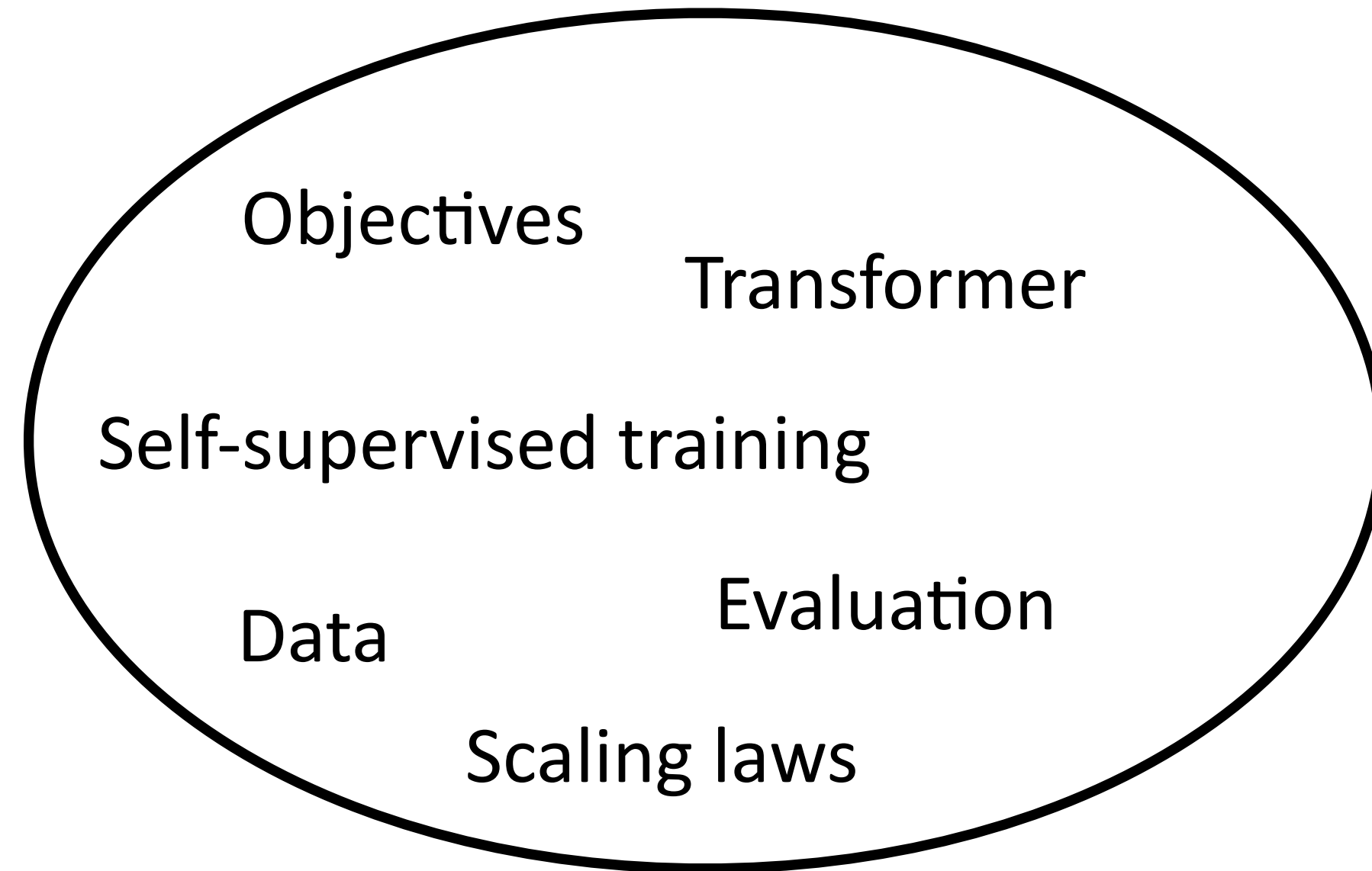


RL

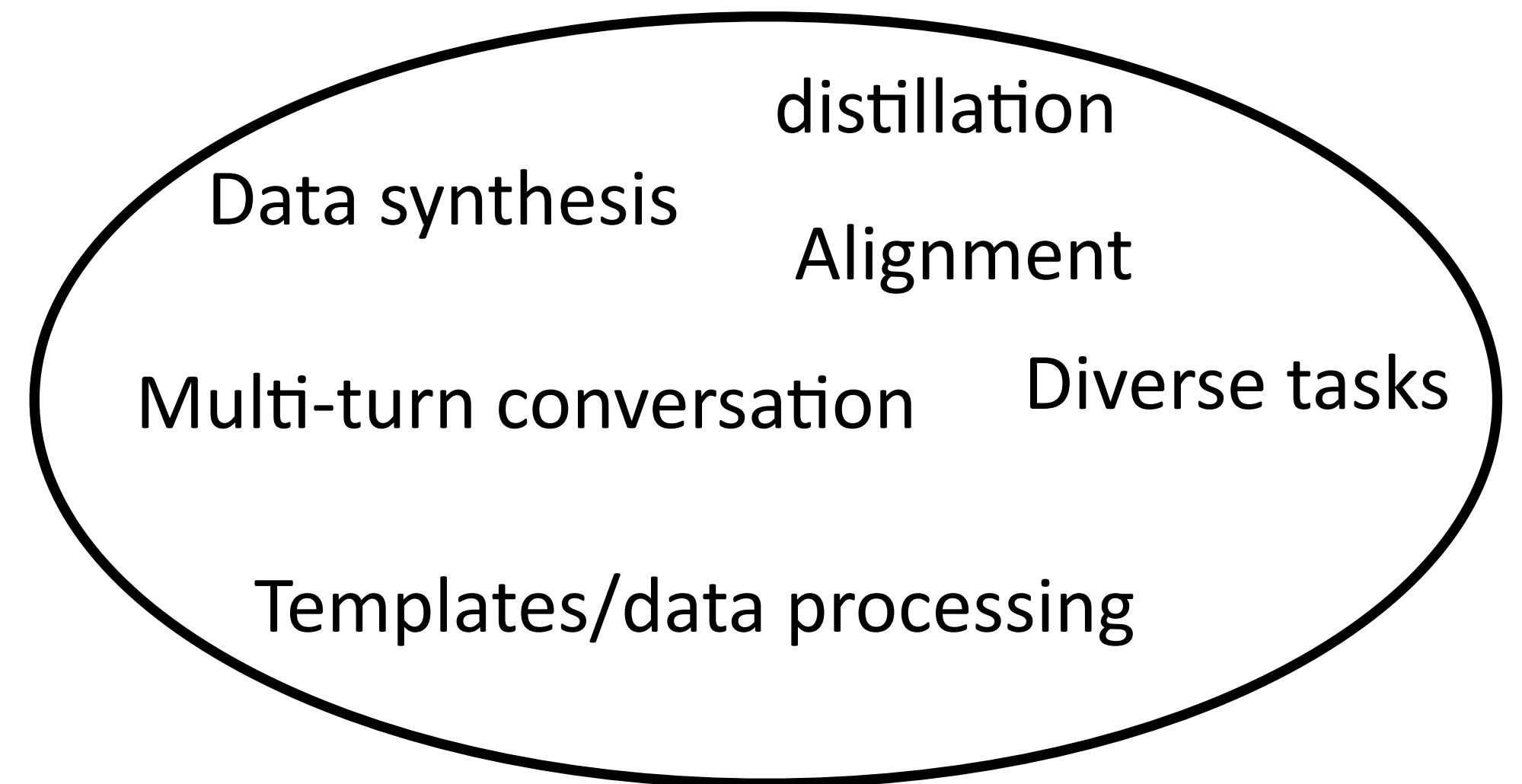


Review — Method

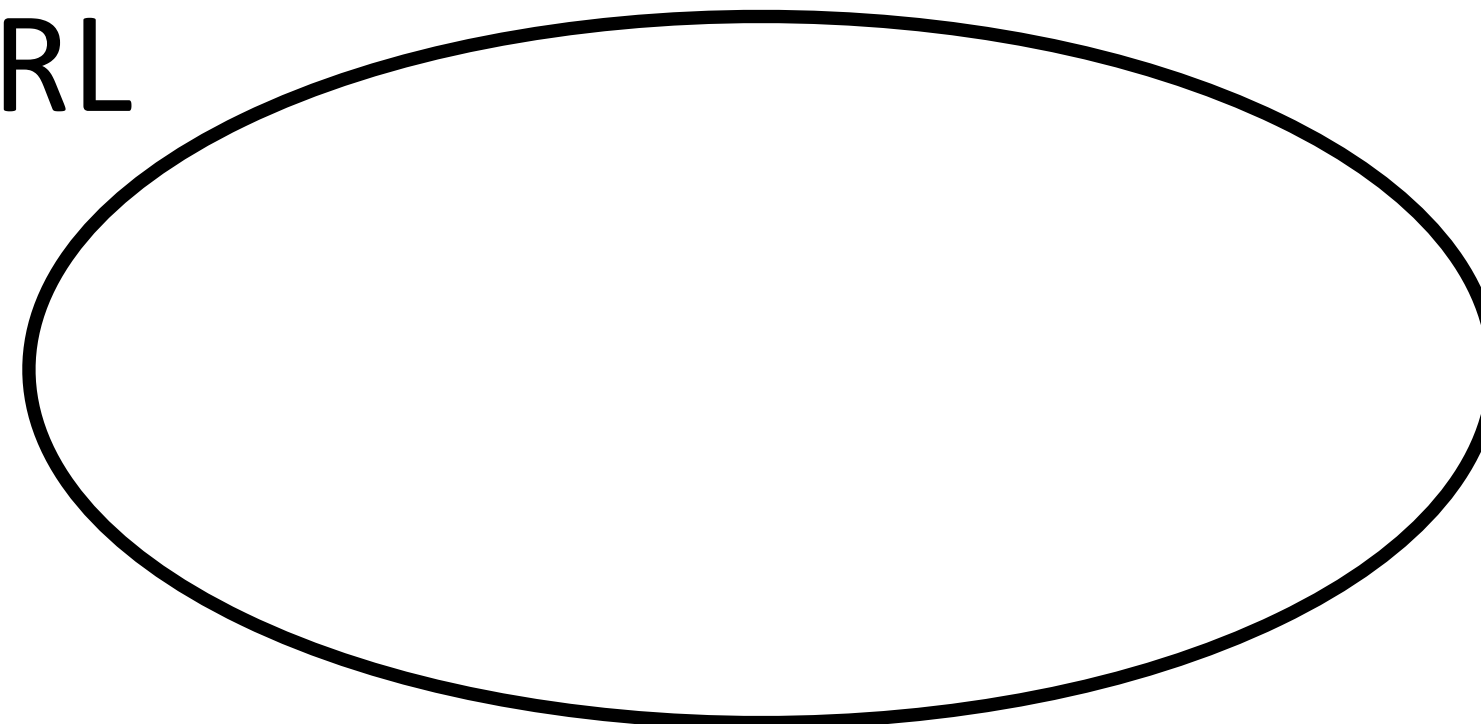
Pretraining



SFT

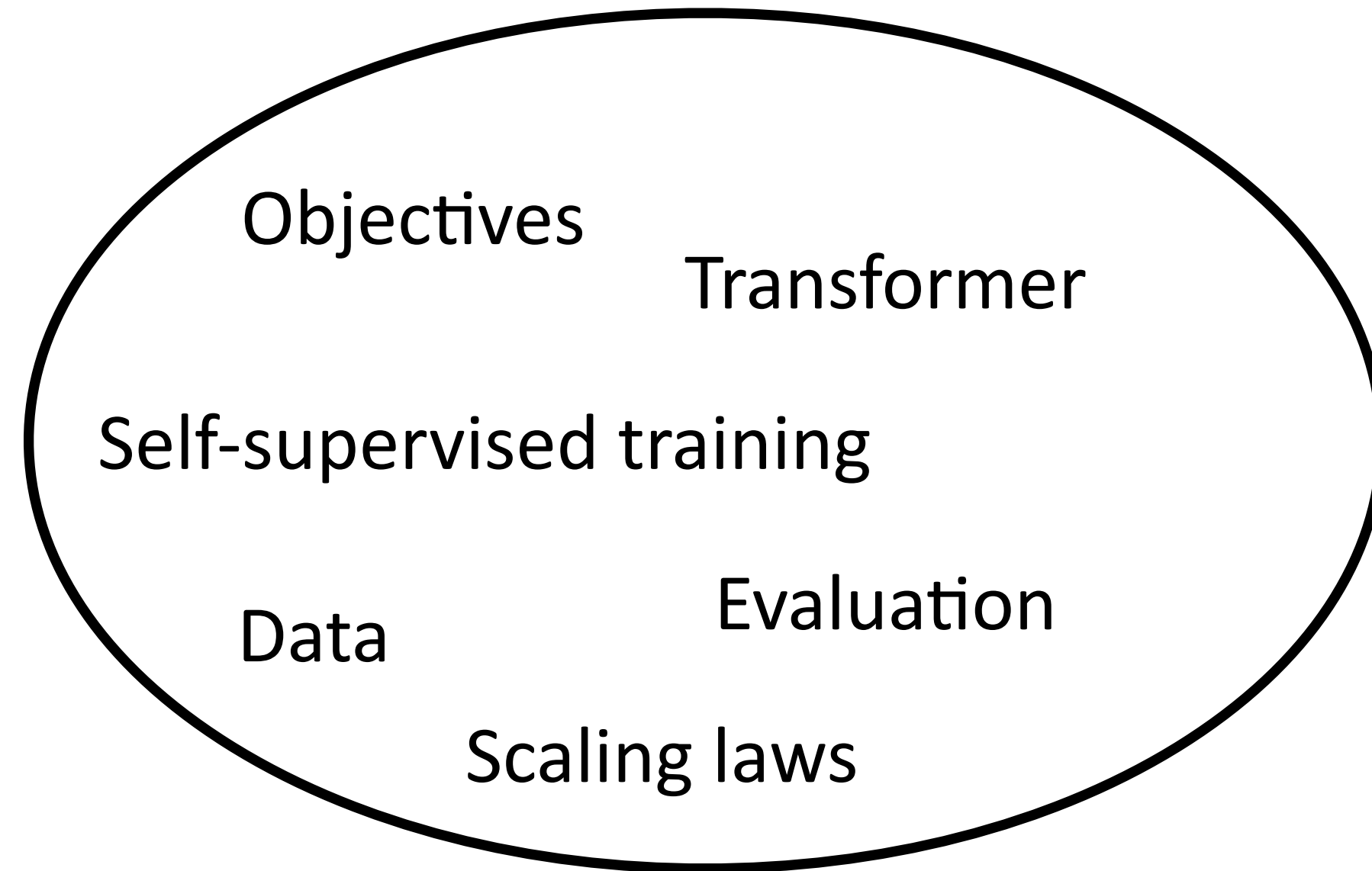


RL

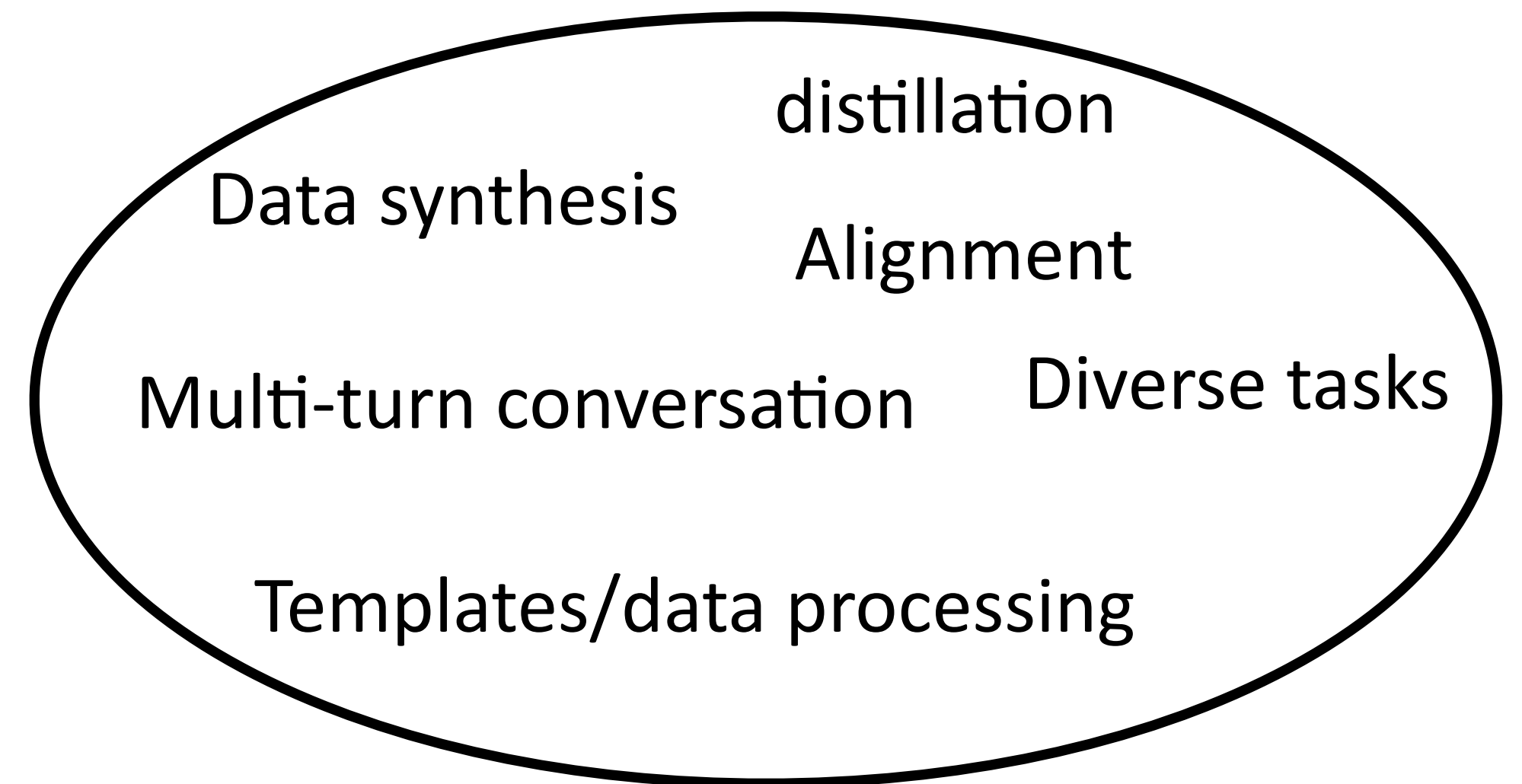


Review — Method

Pretraining



SFT

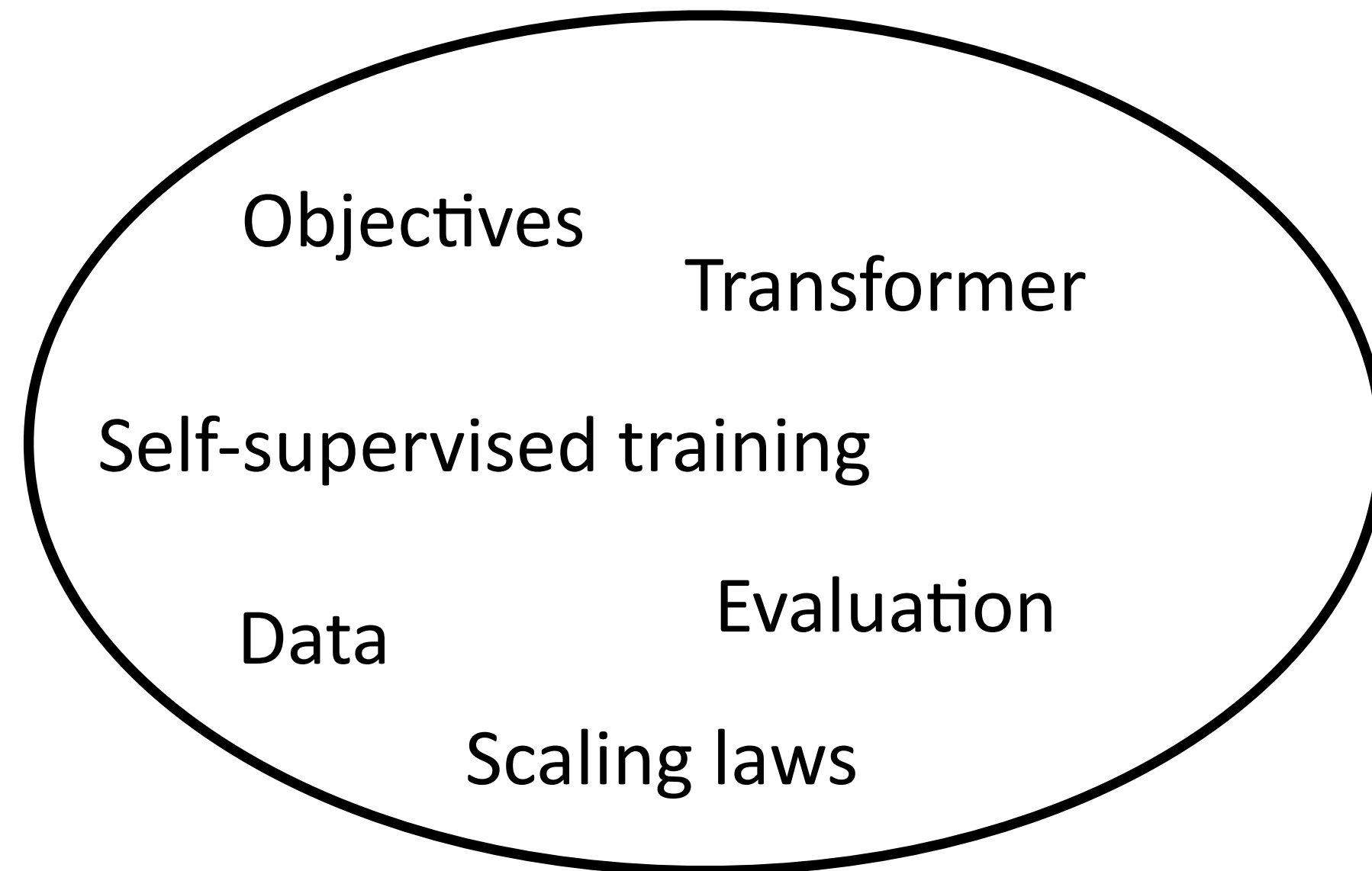


RL

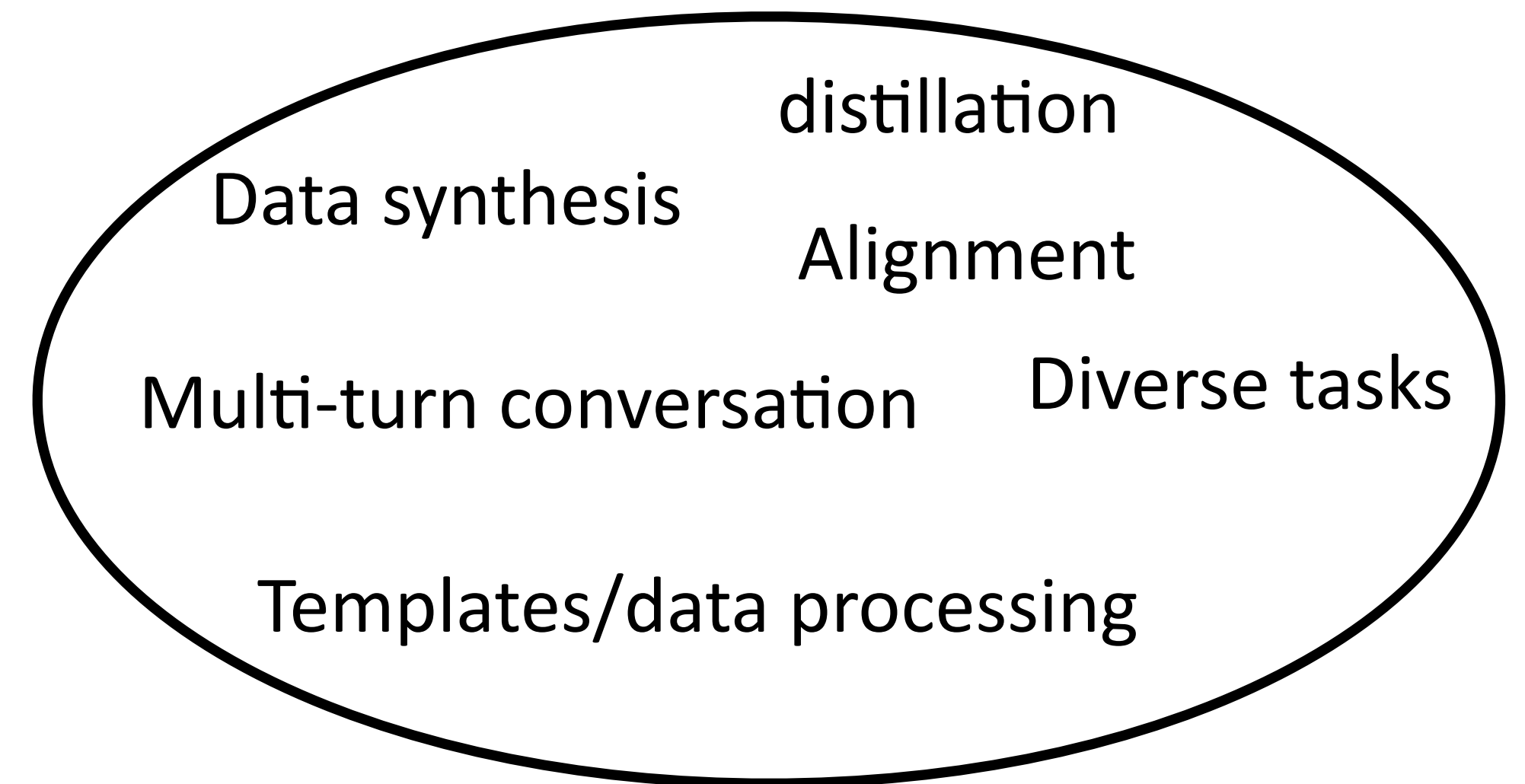


Review — Method

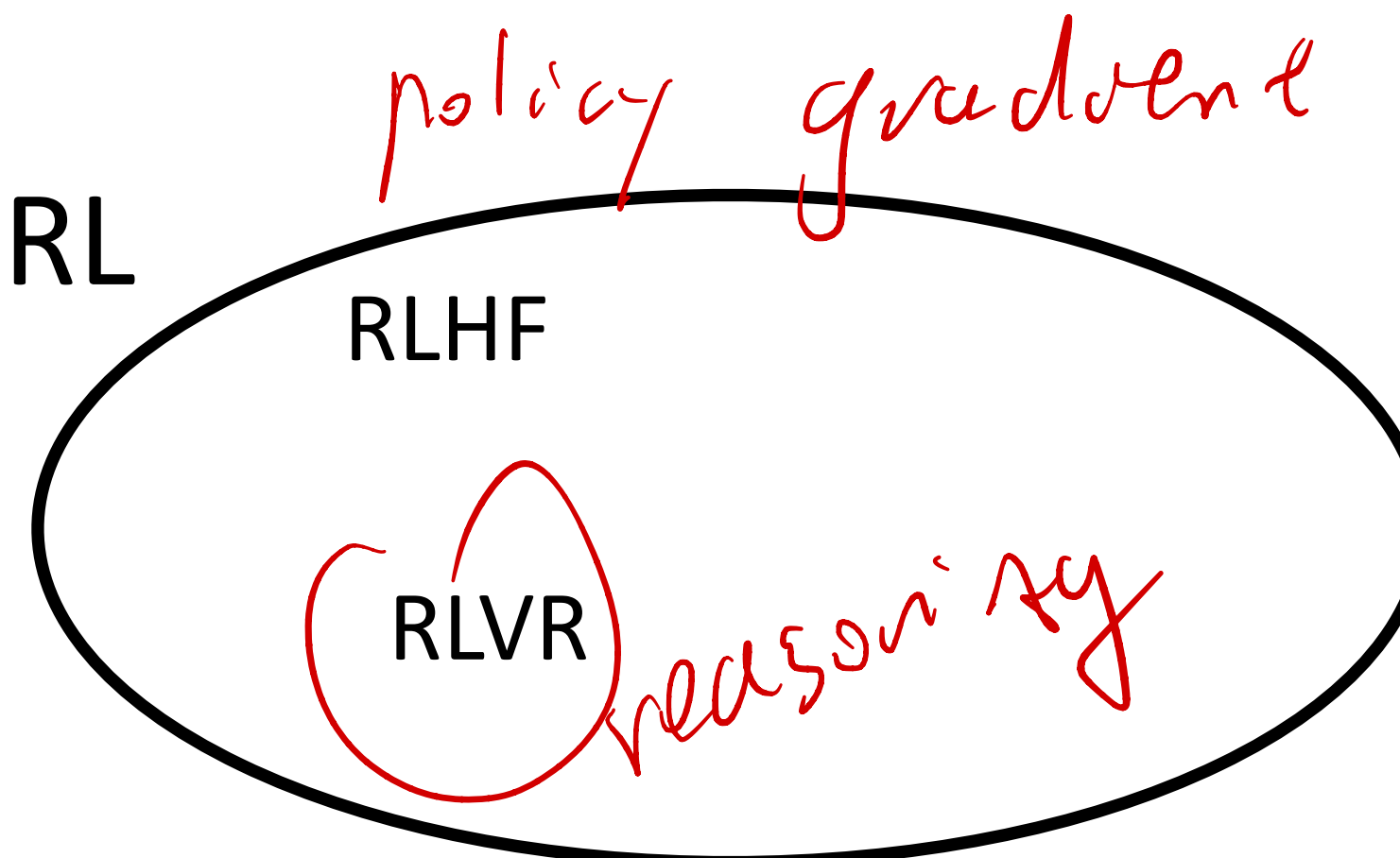
Pretraining



SFT

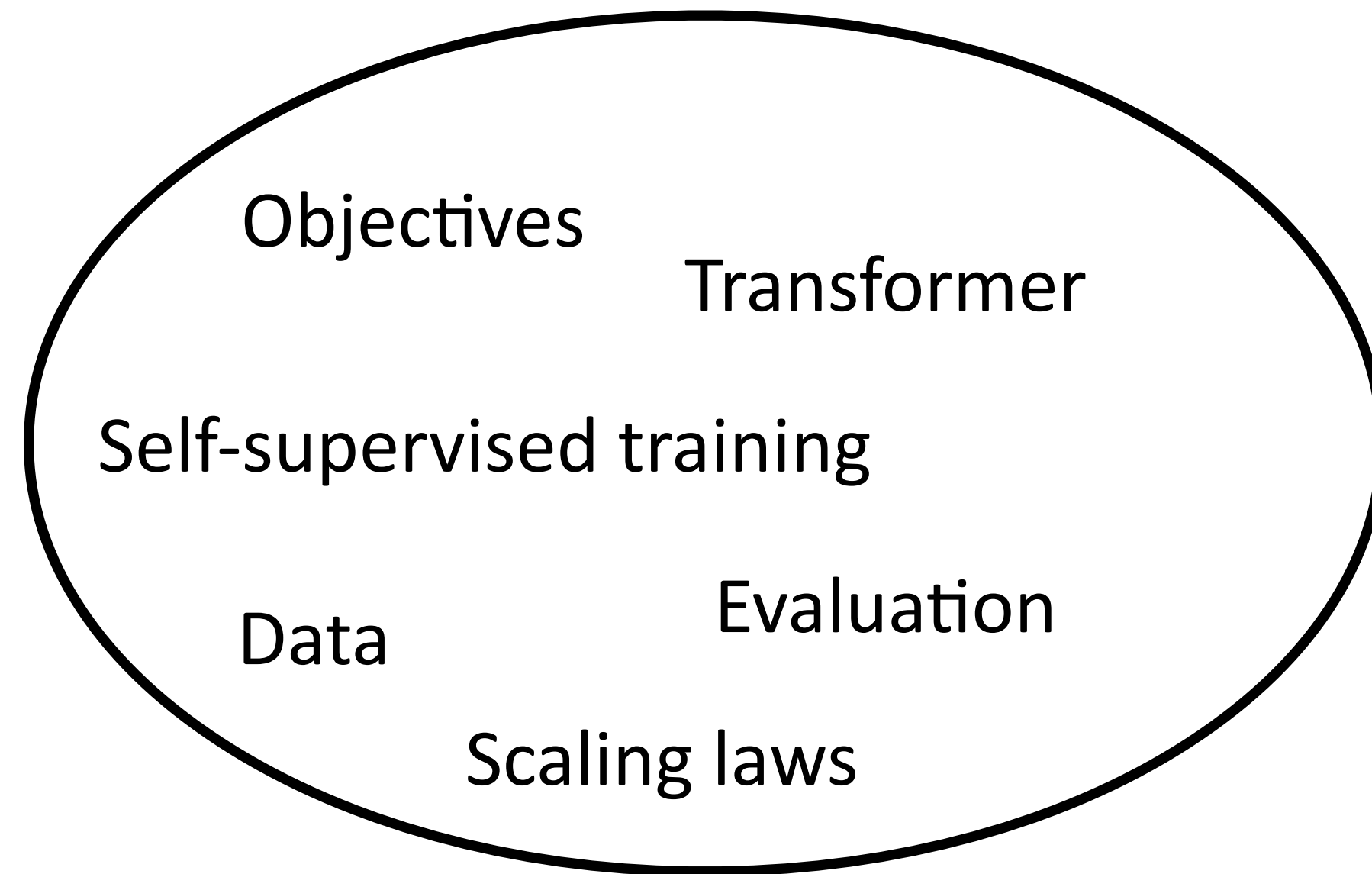


RL

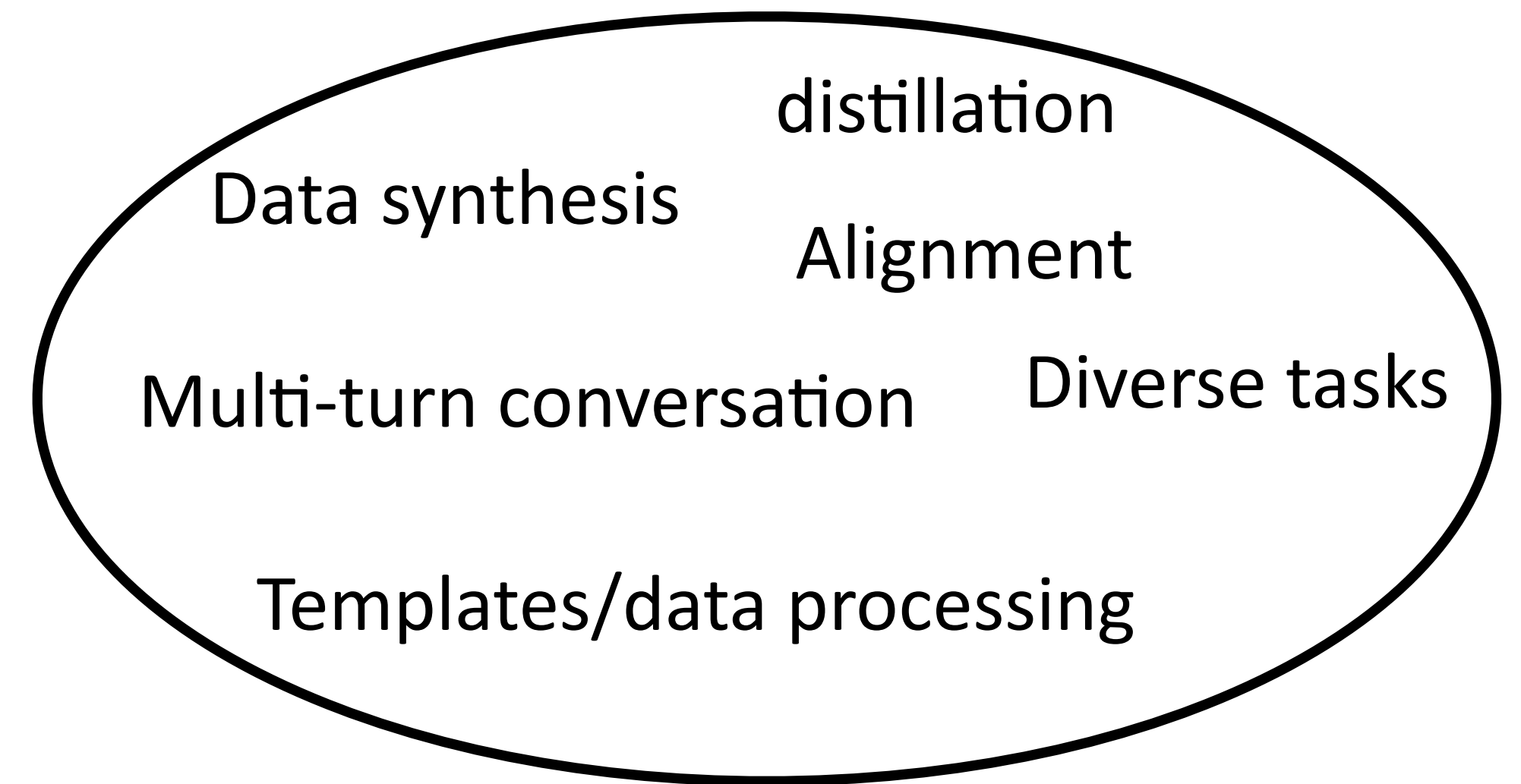


Review — Method

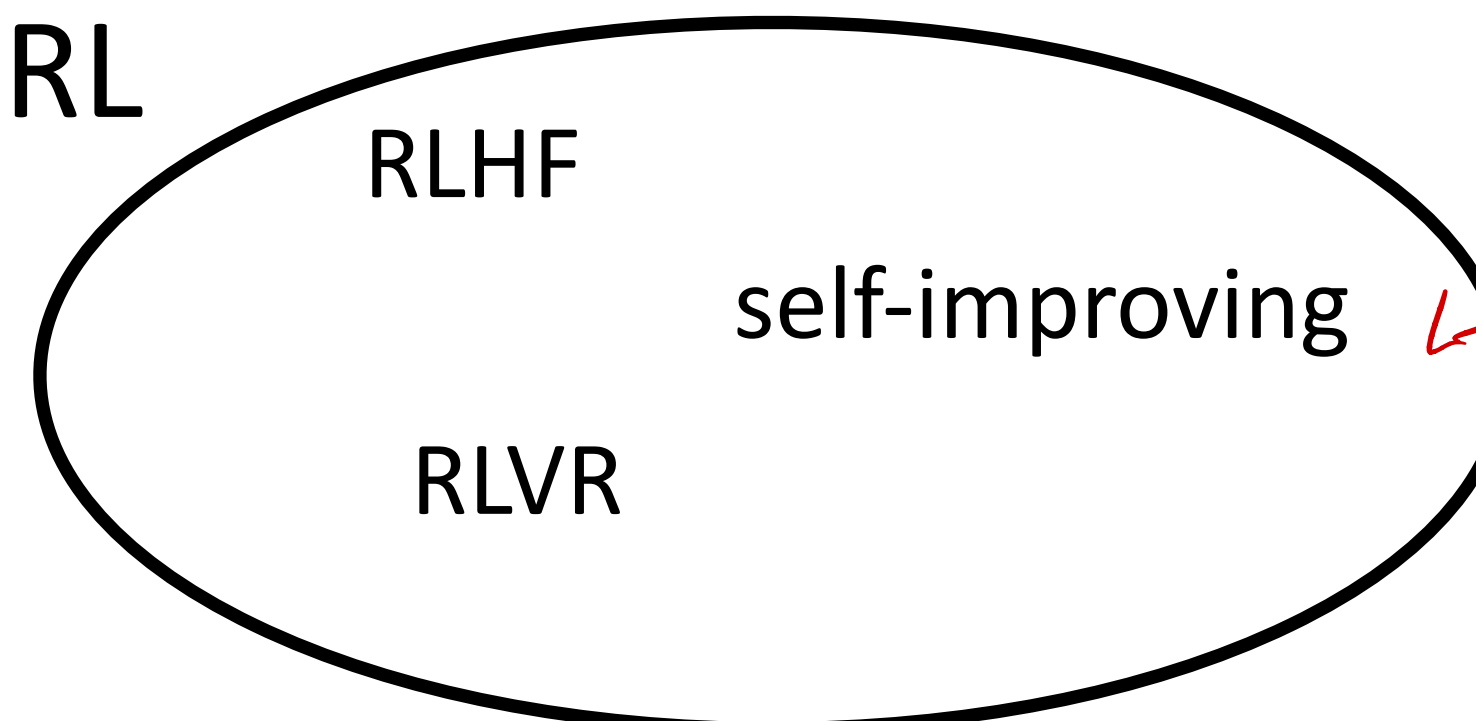
Pretraining



SFT



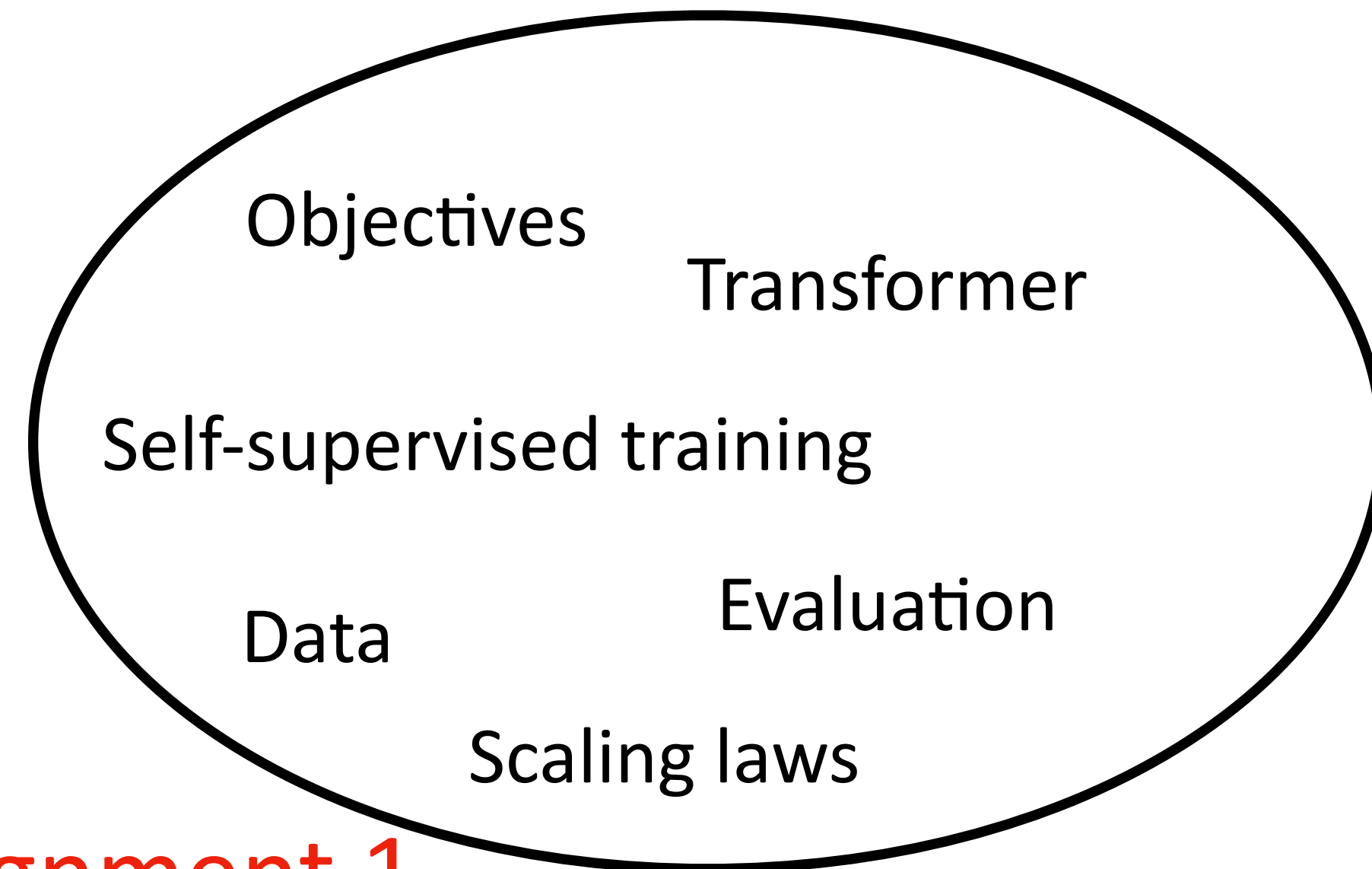
RL



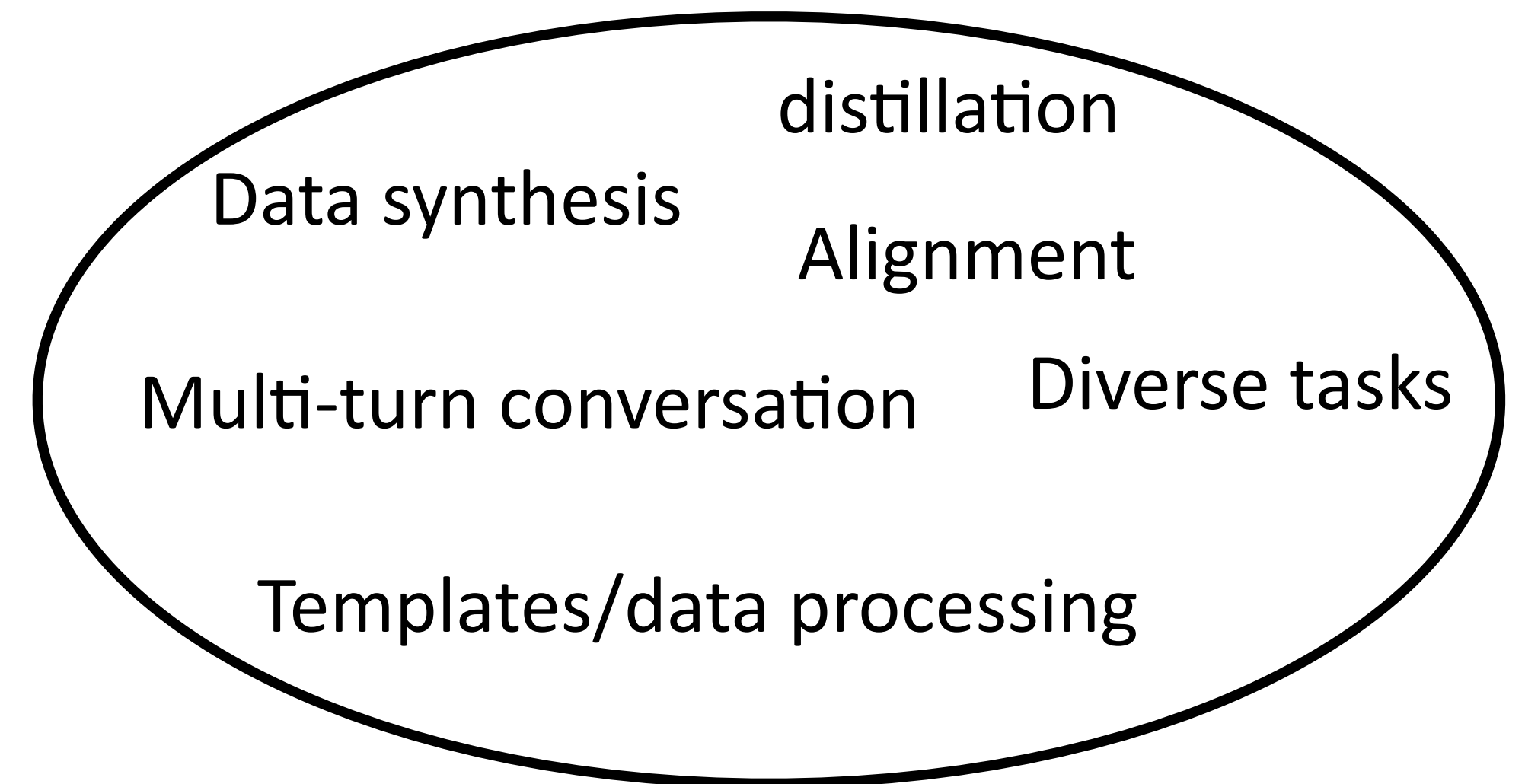
no imitation learning

Review — Method

Pretraining

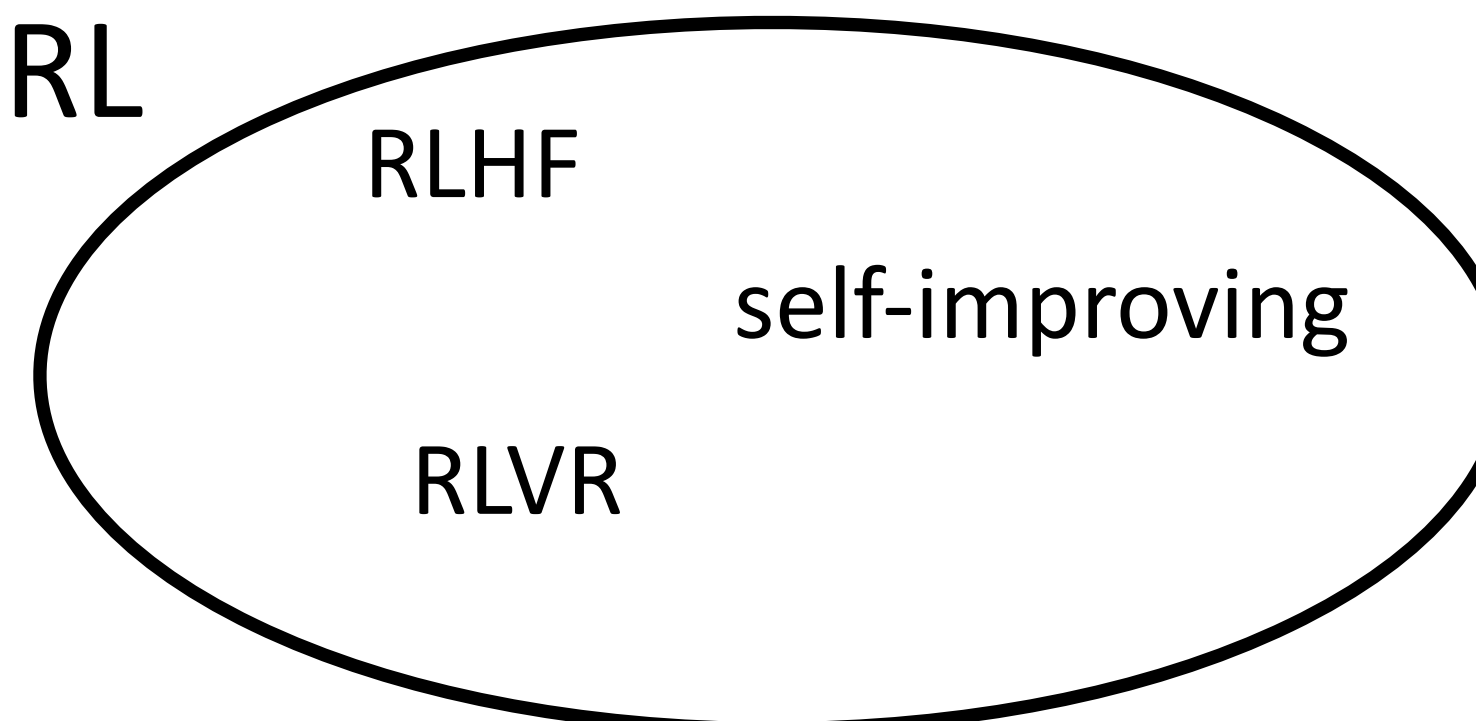


SFT



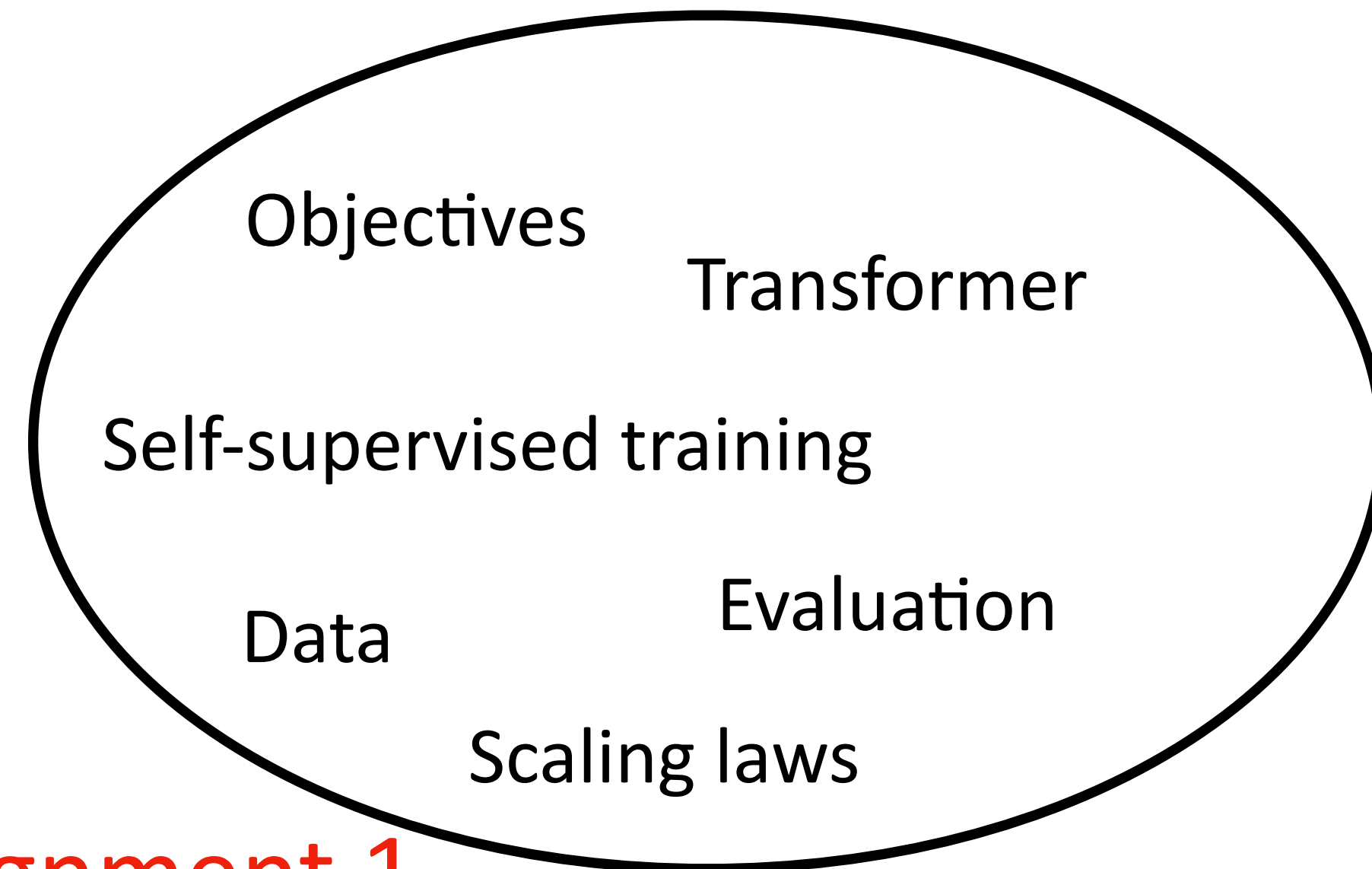
Assignment 1

RL



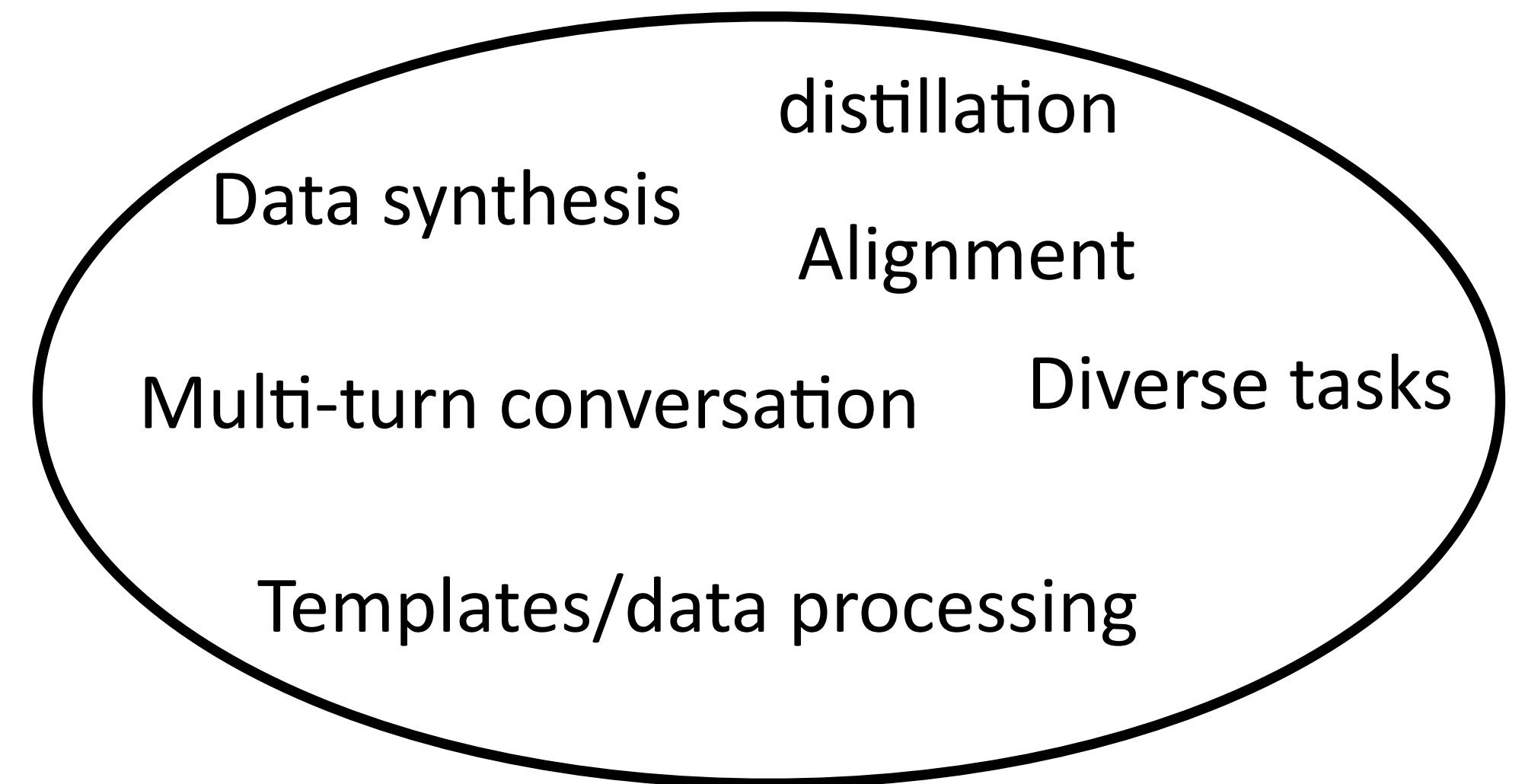
Review — Method

Pretraining



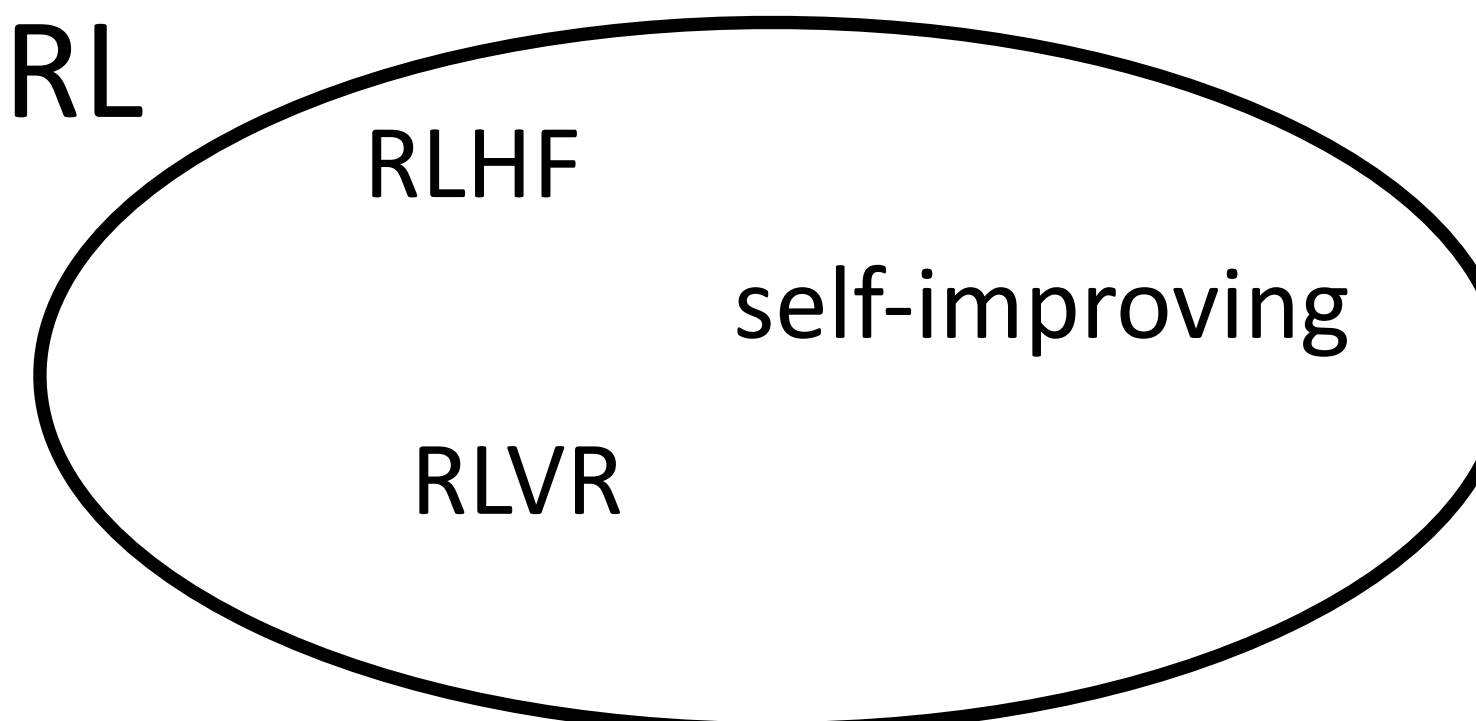
Assignment 1

SFT



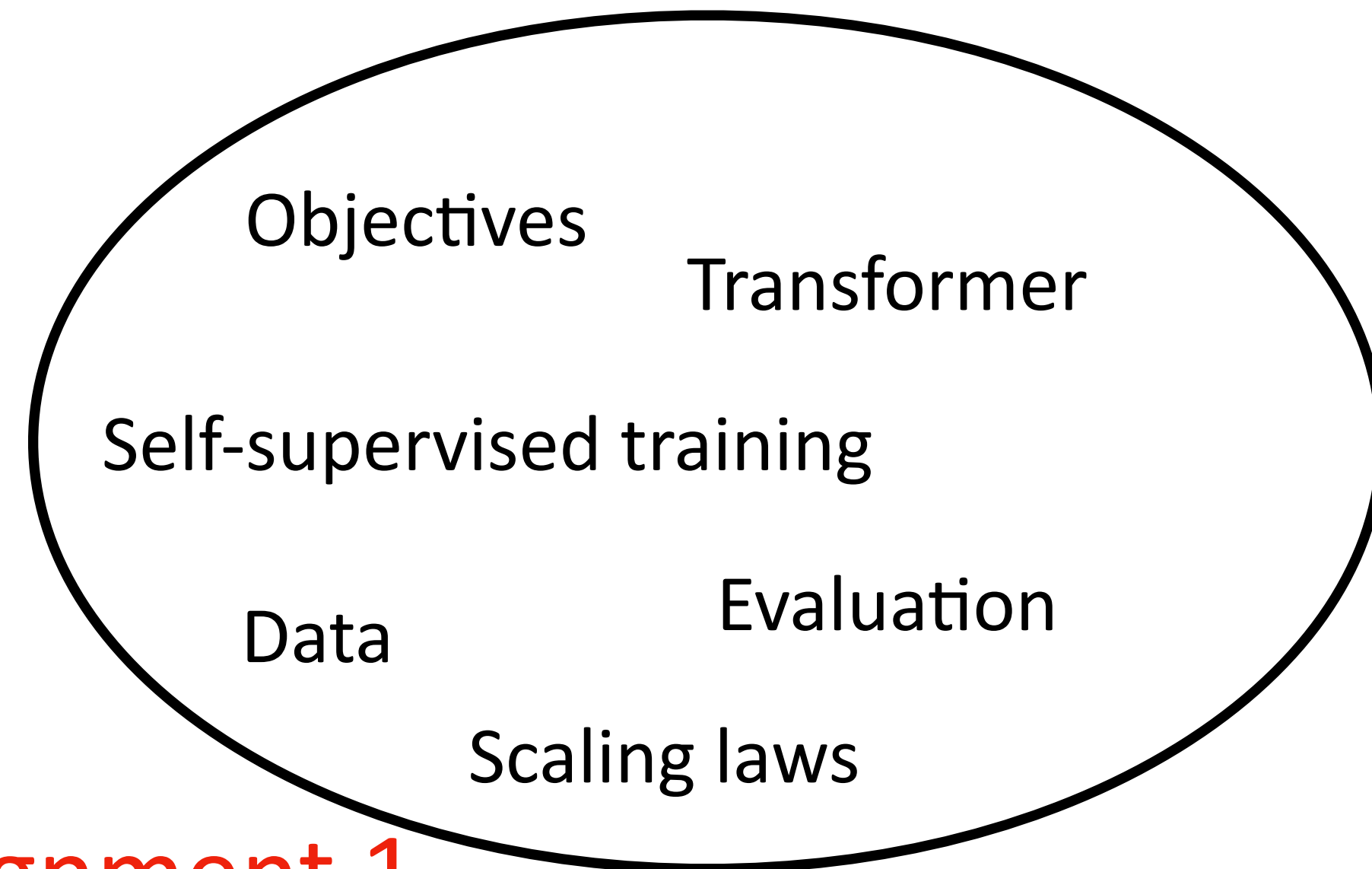
Assignment 2

RL



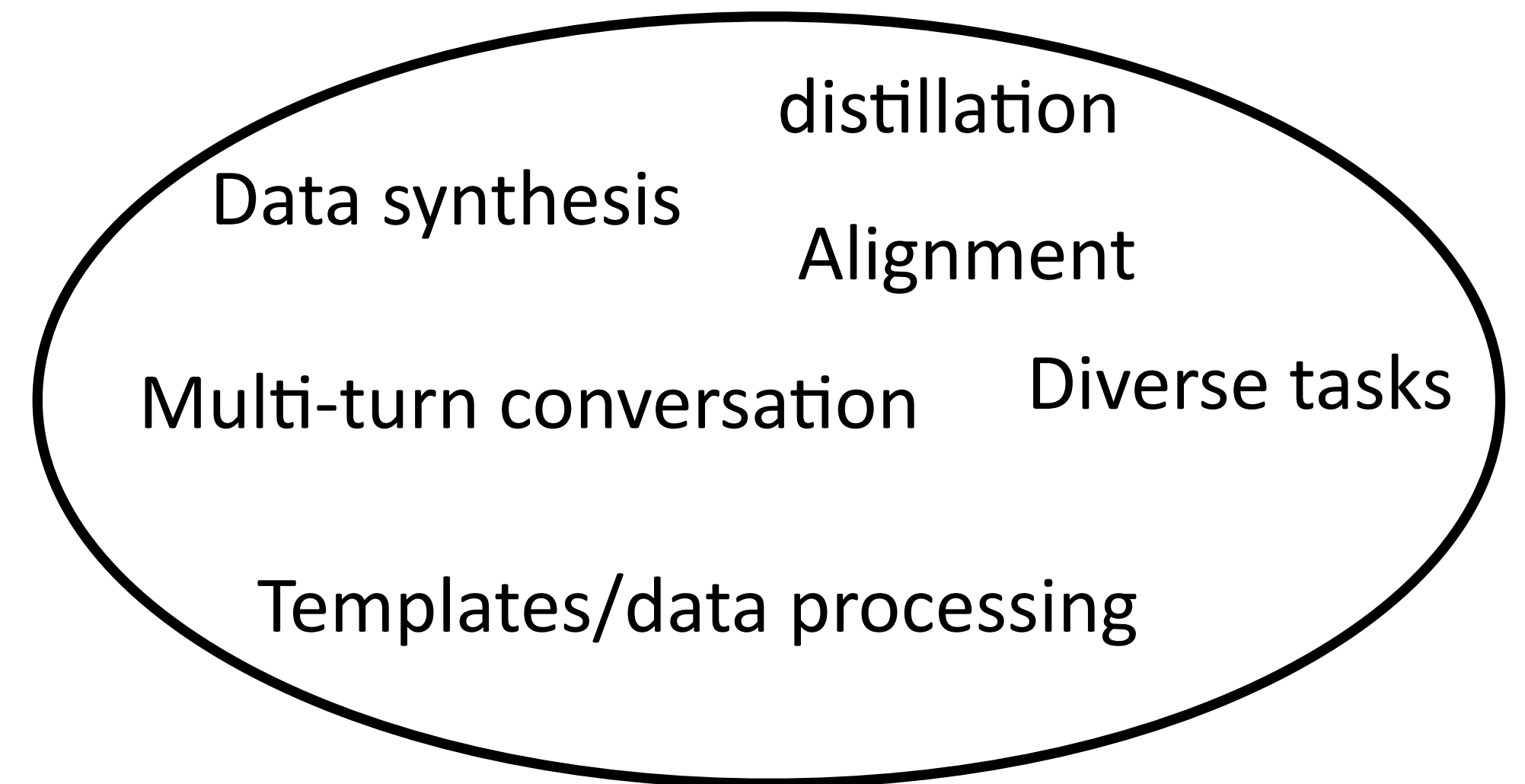
Review — Method

Pretraining



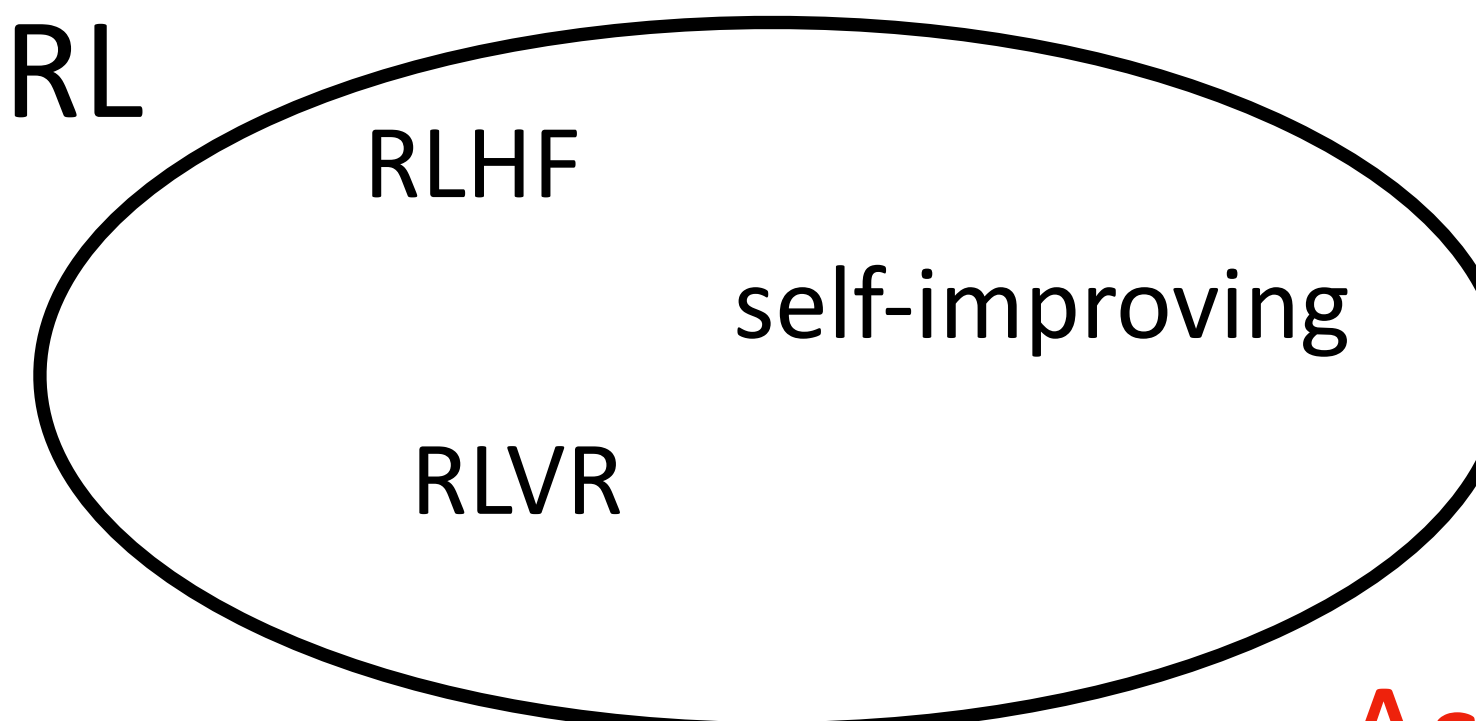
Assignment 1

SFT



Assignment 2

RL

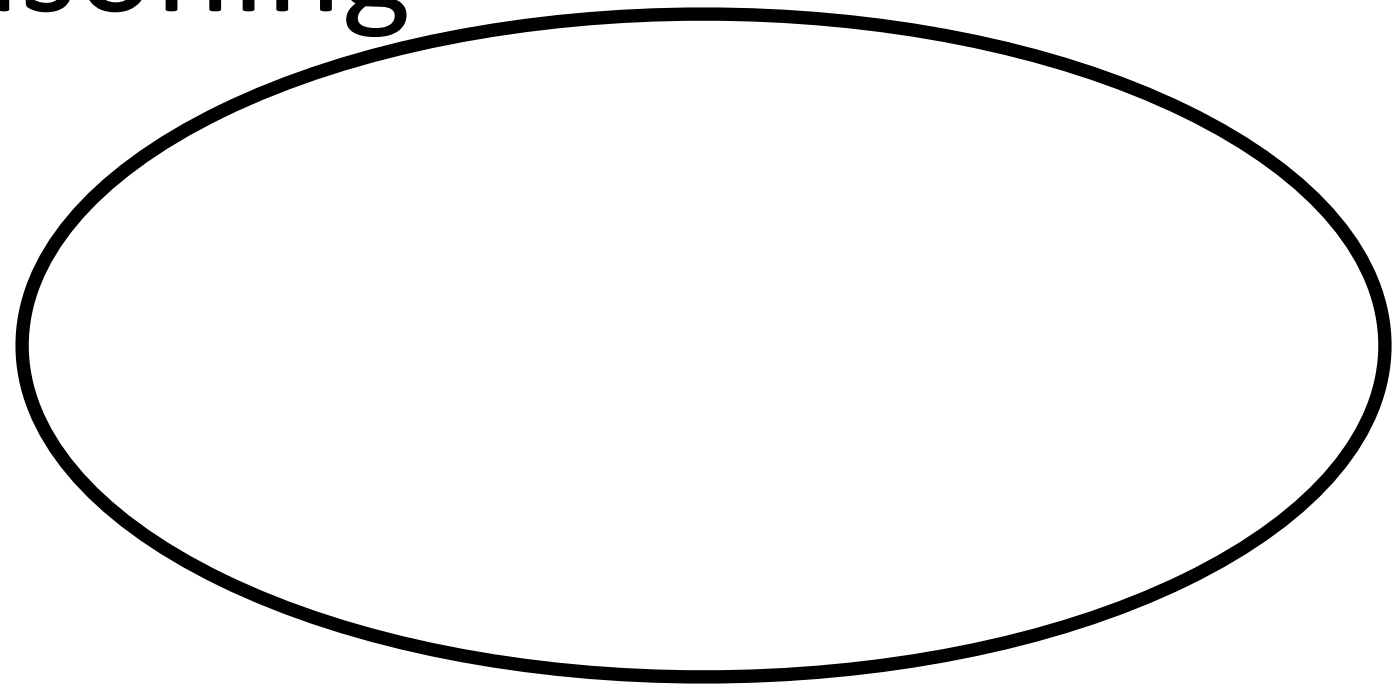


Assignment 3

Review — Applications

Review — Applications

Reasoning



Review — Applications

Reasoning



chain-of-thought

Review — Applications

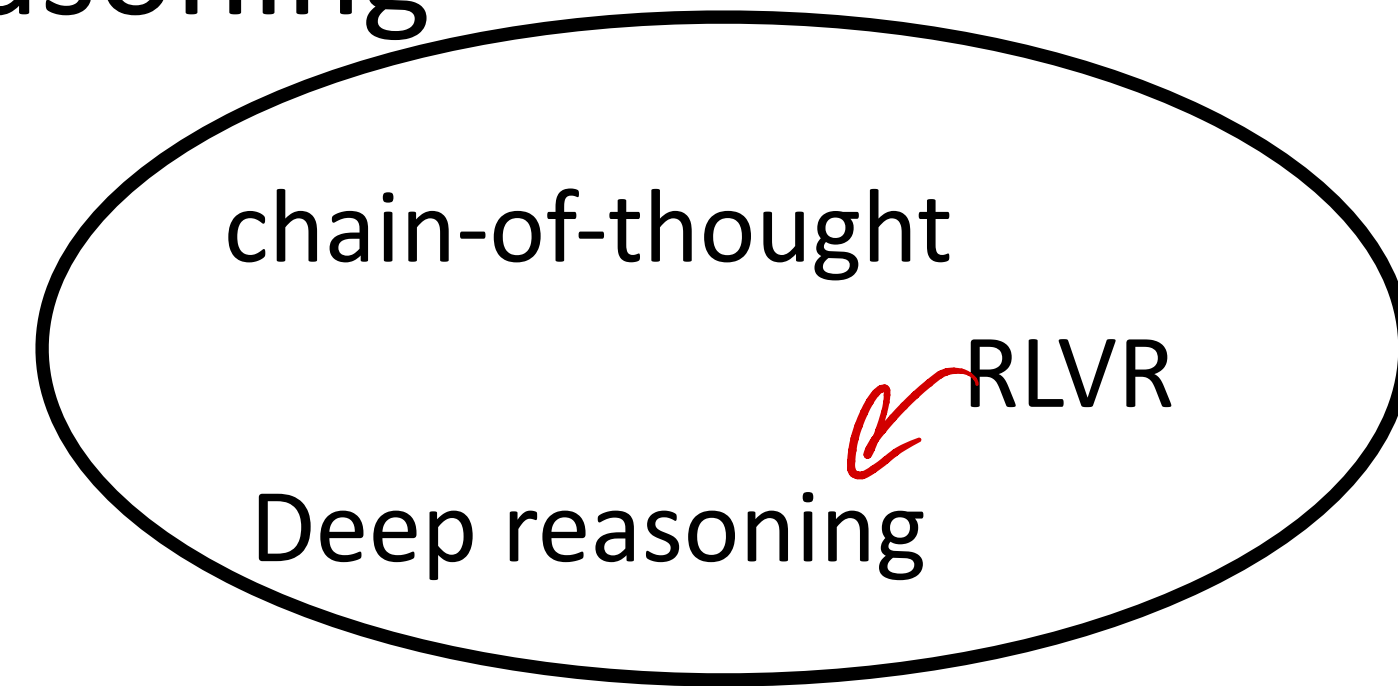
Reasoning

chain-of-thought

Deep reasoning

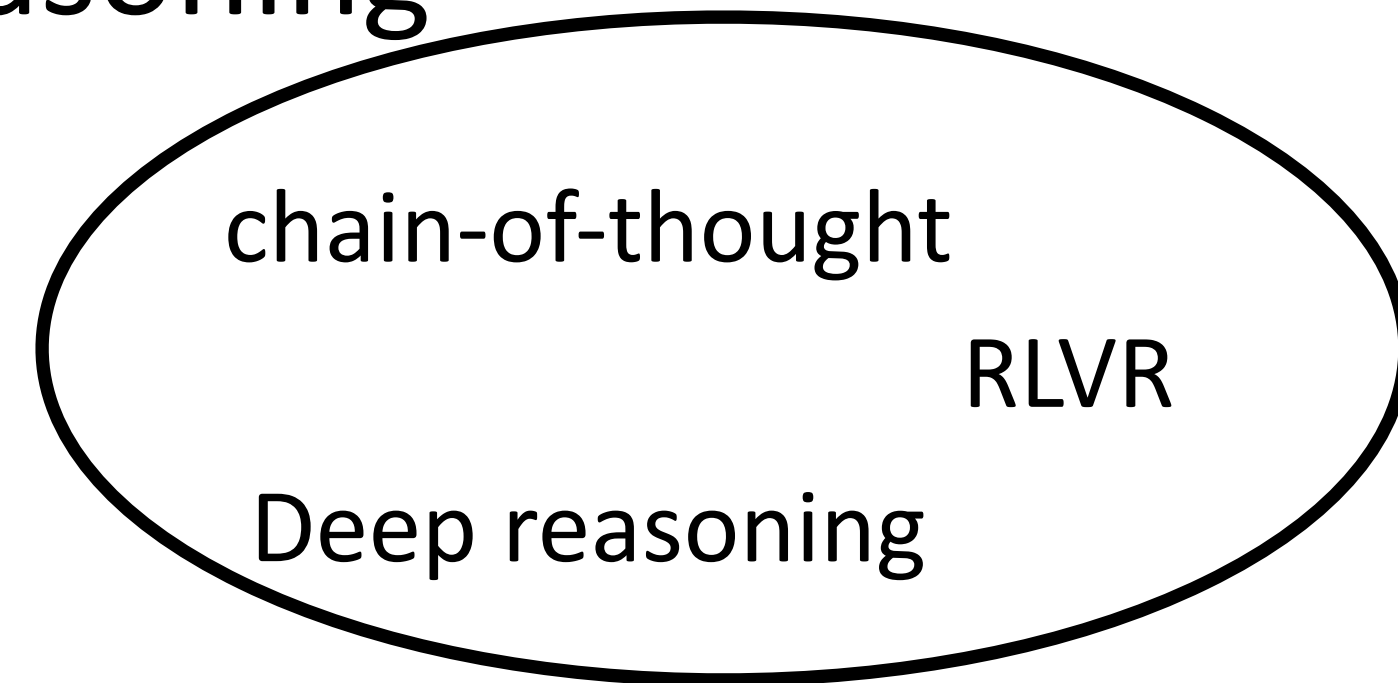
Review — Applications

Reasoning



Review — Applications

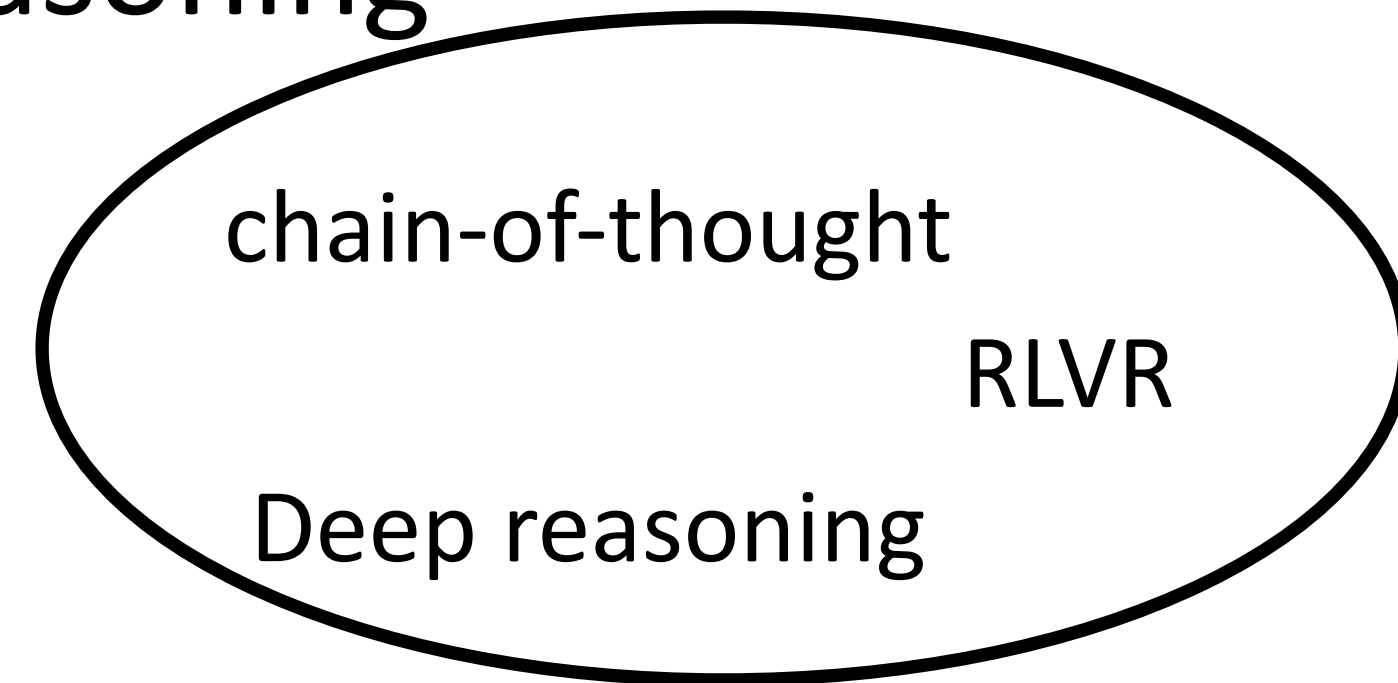
Reasoning



Assignment 3, 4

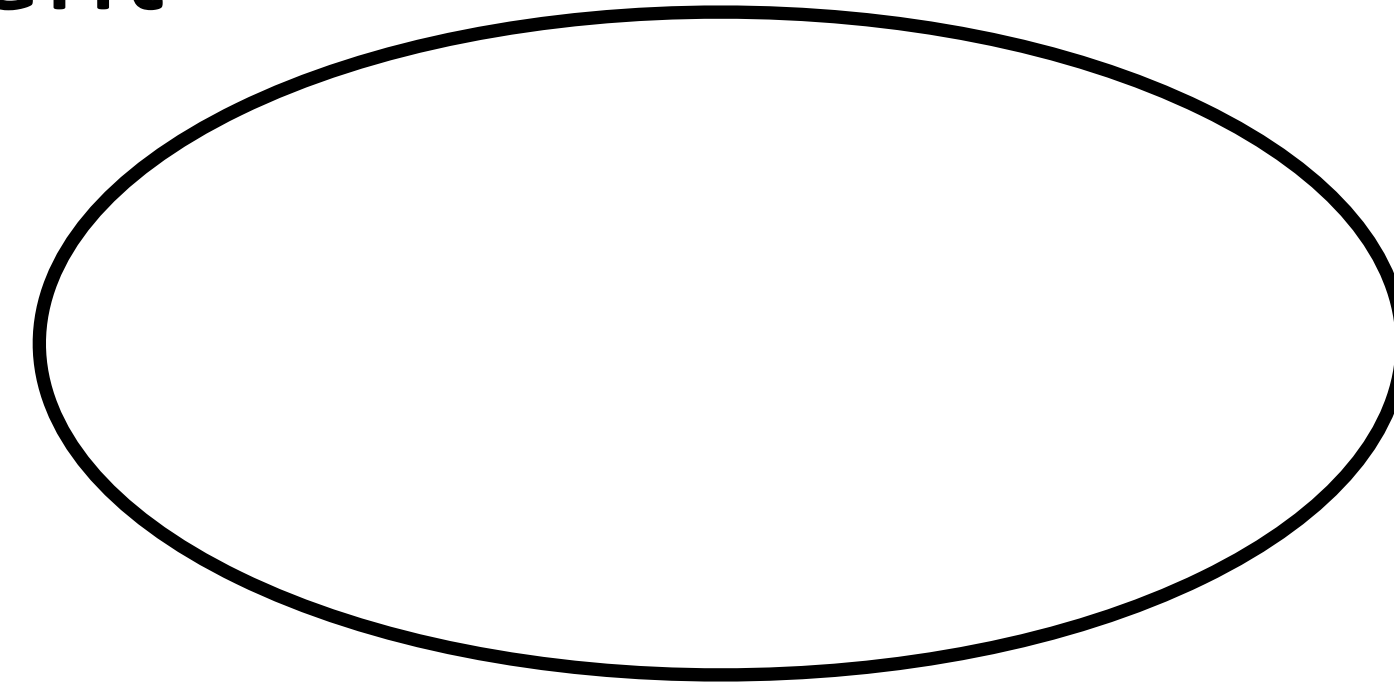
Review — Applications

Reasoning



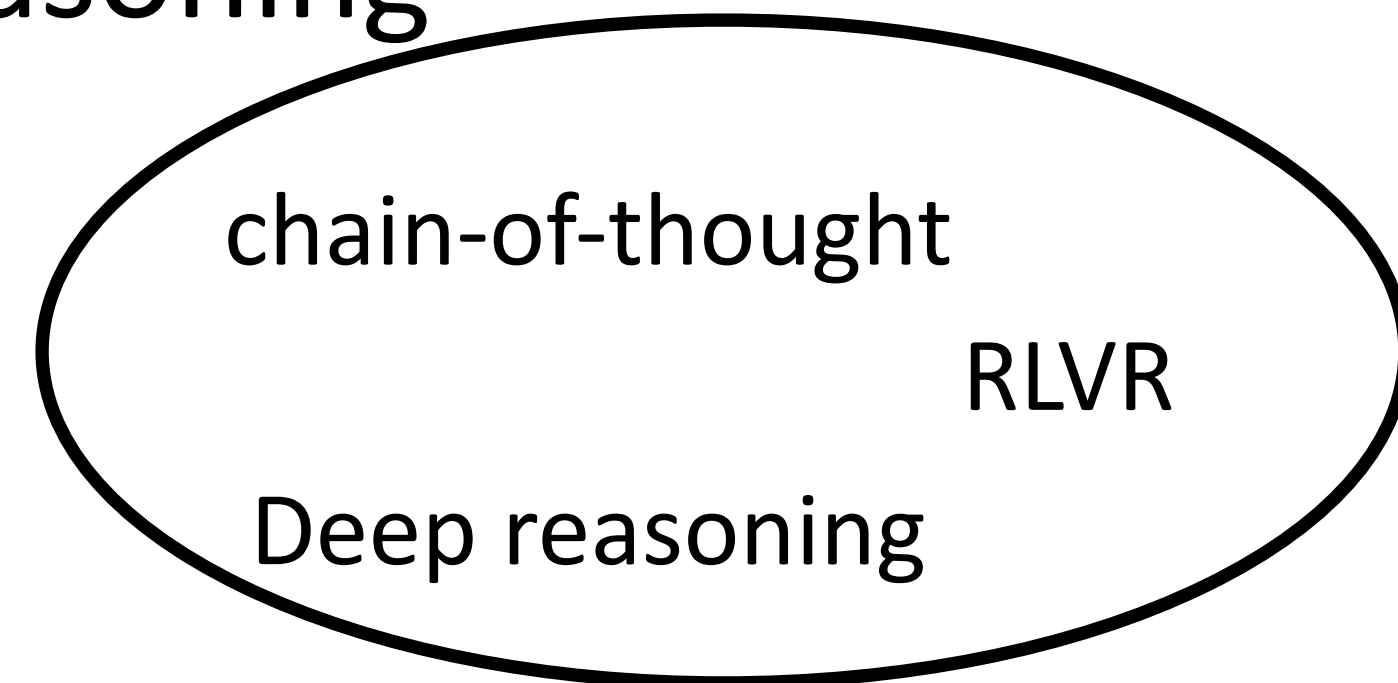
Assignment 3, 4

Agent



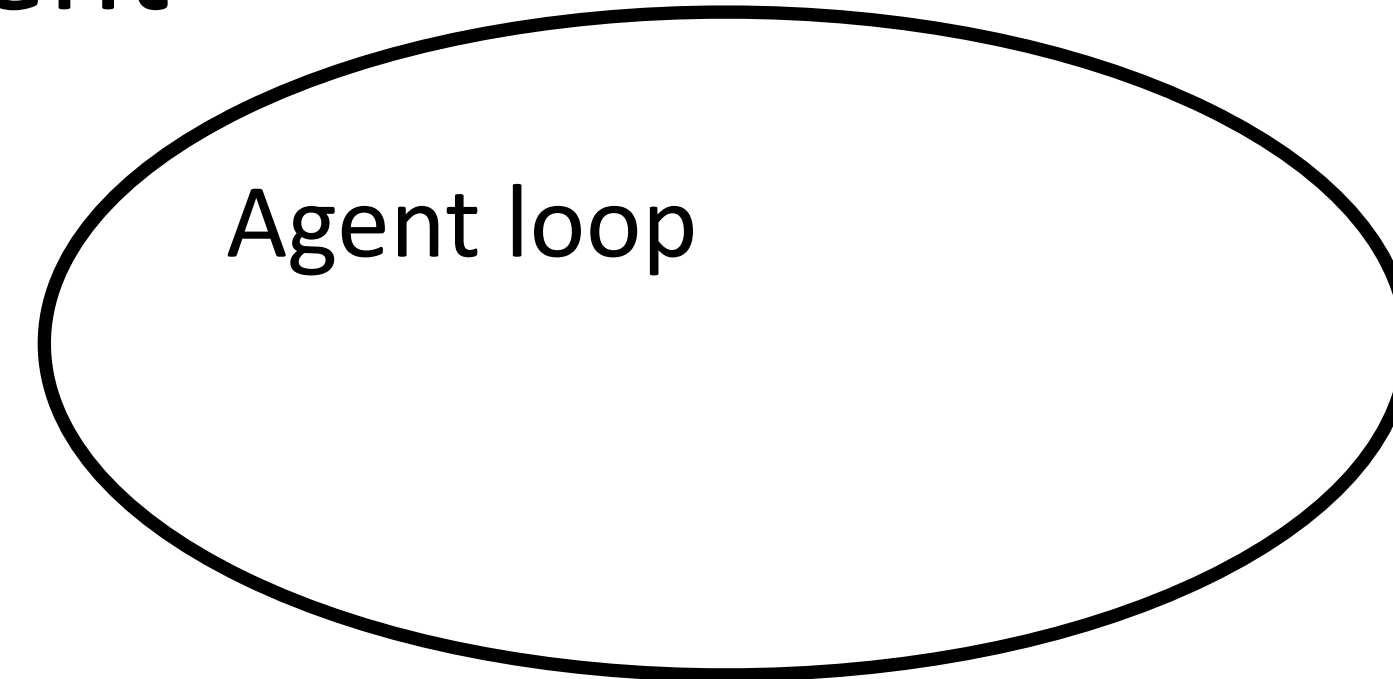
Review — Applications

Reasoning



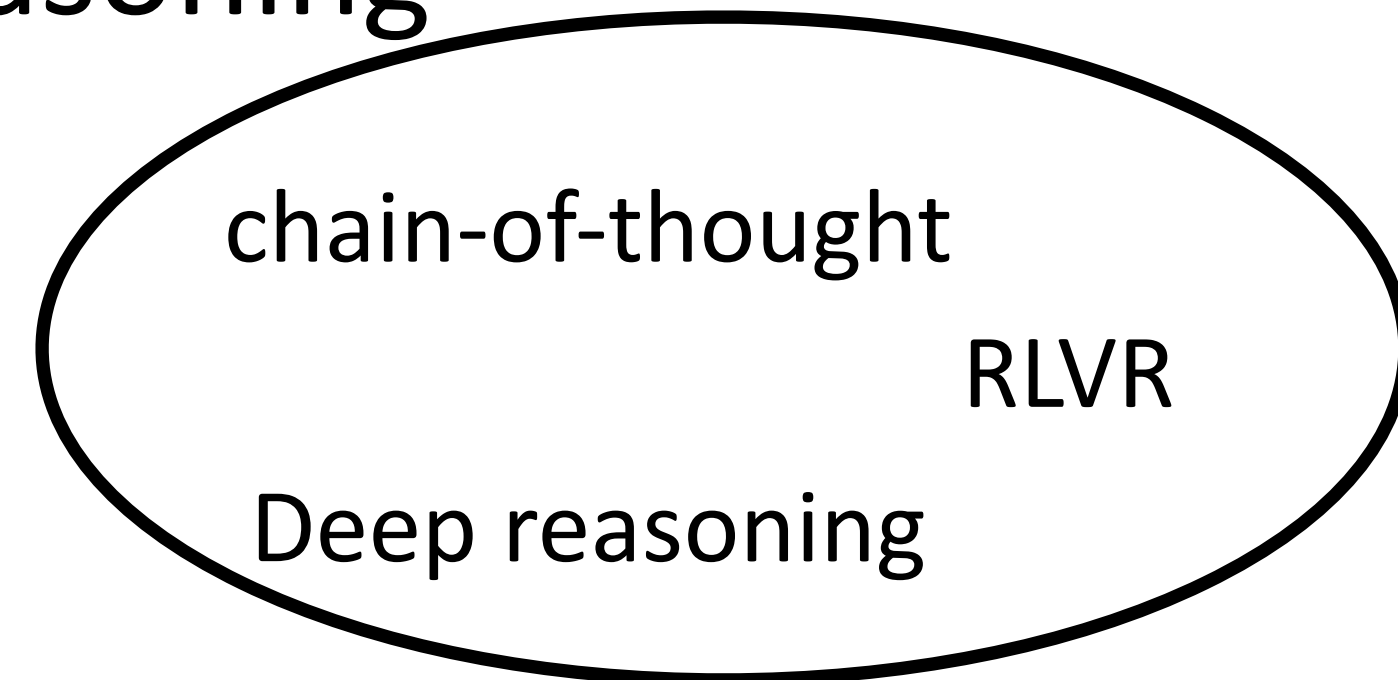
Assignment 3, 4

Agent



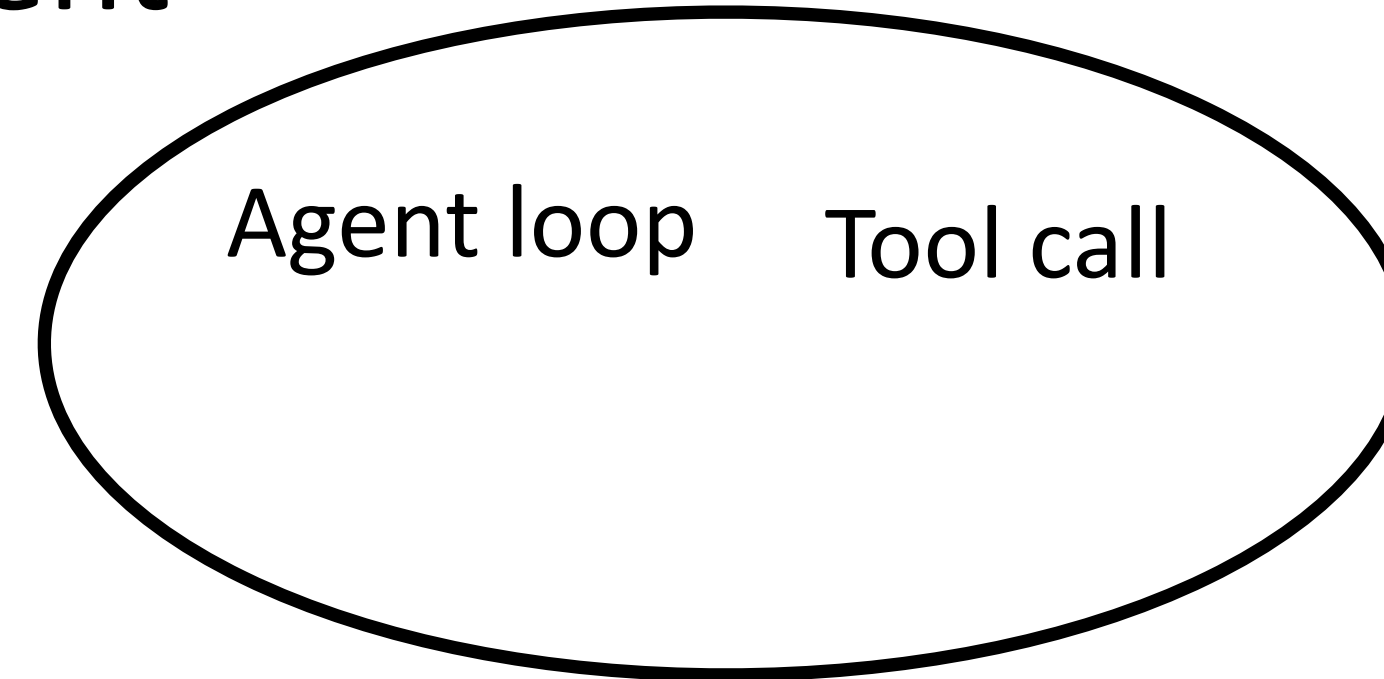
Review — Applications

Reasoning



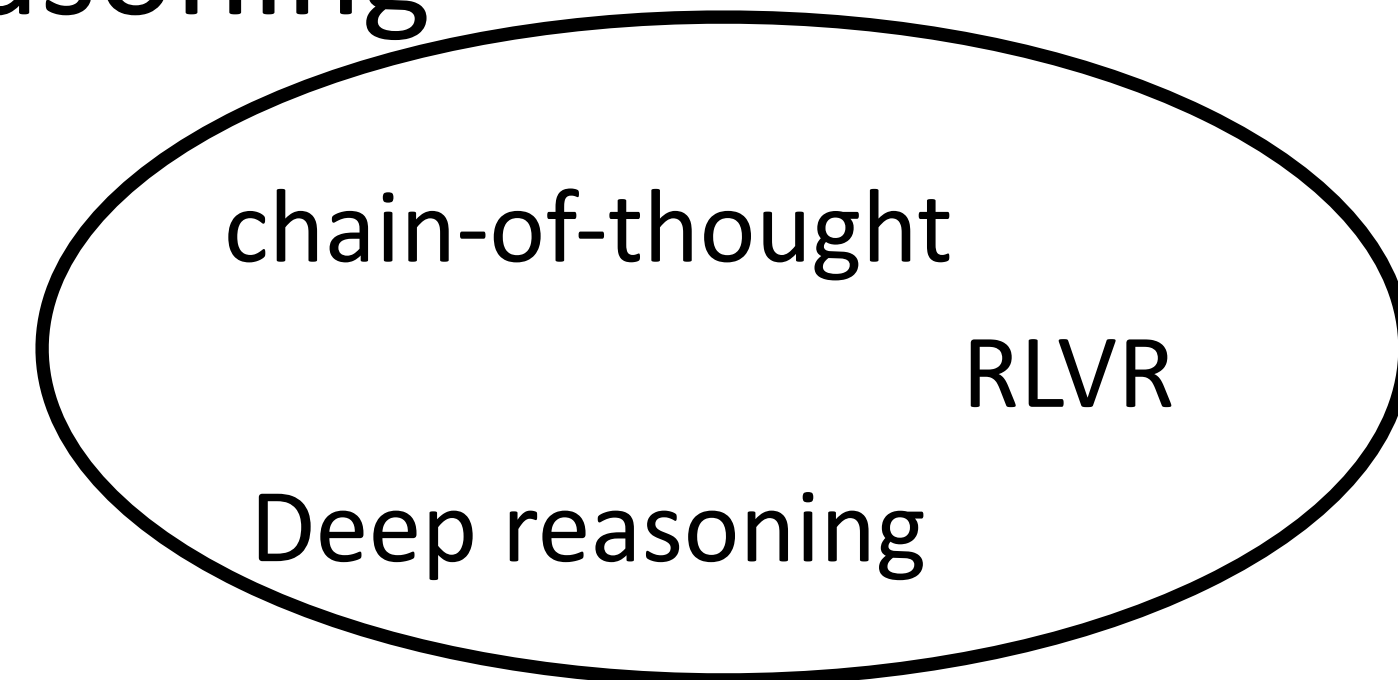
Assignment 3, 4

Agent



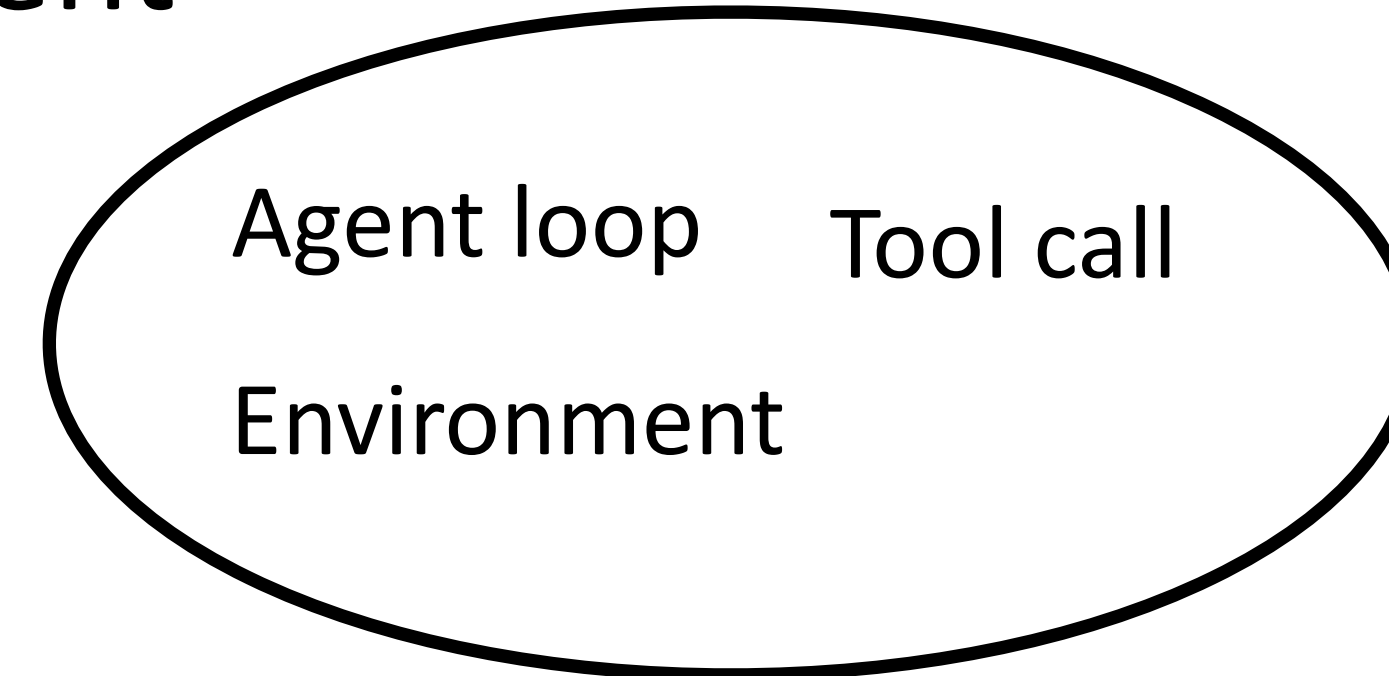
Review — Applications

Reasoning



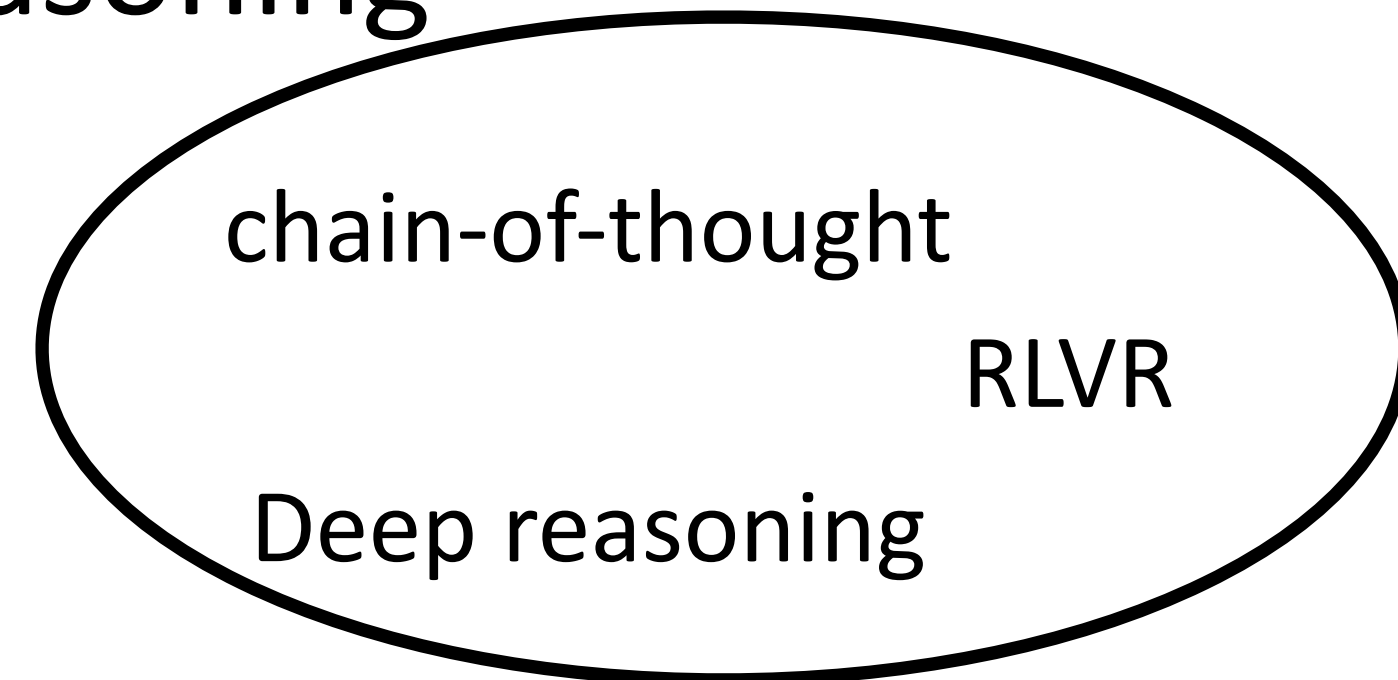
Assignment 3, 4

Agent



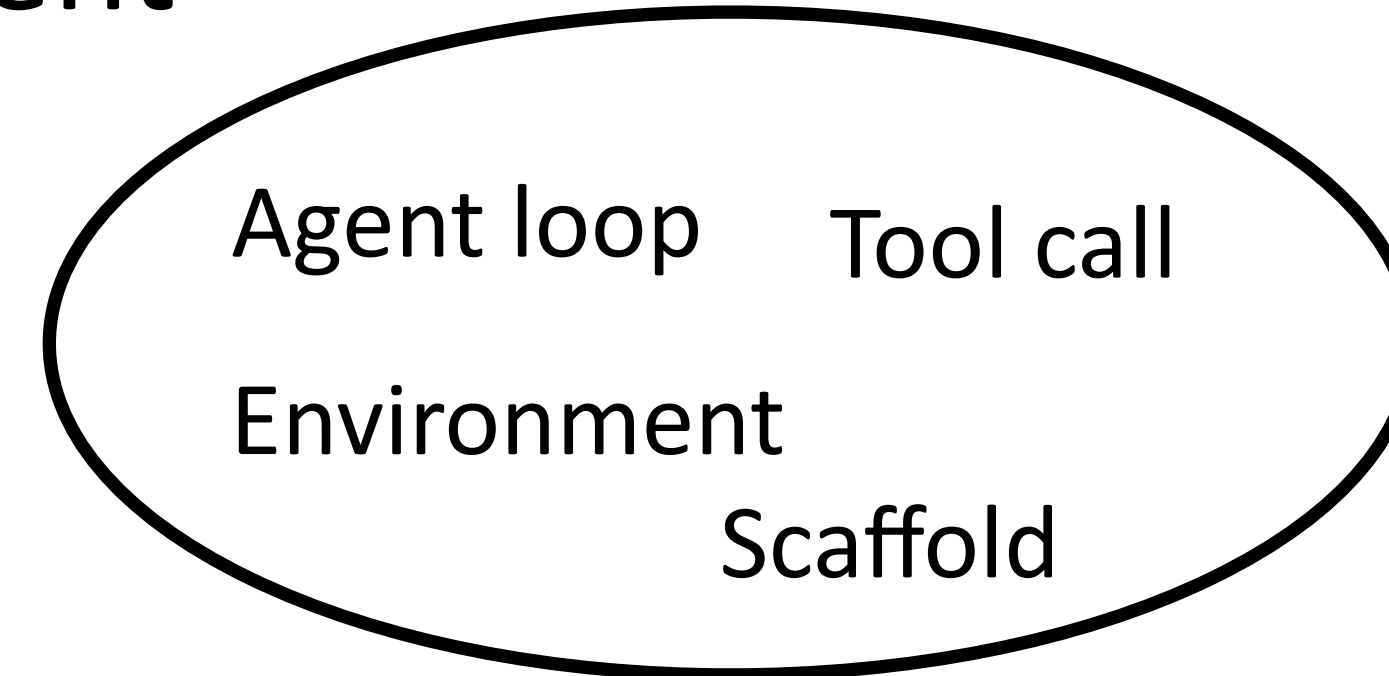
Review — Applications

Reasoning



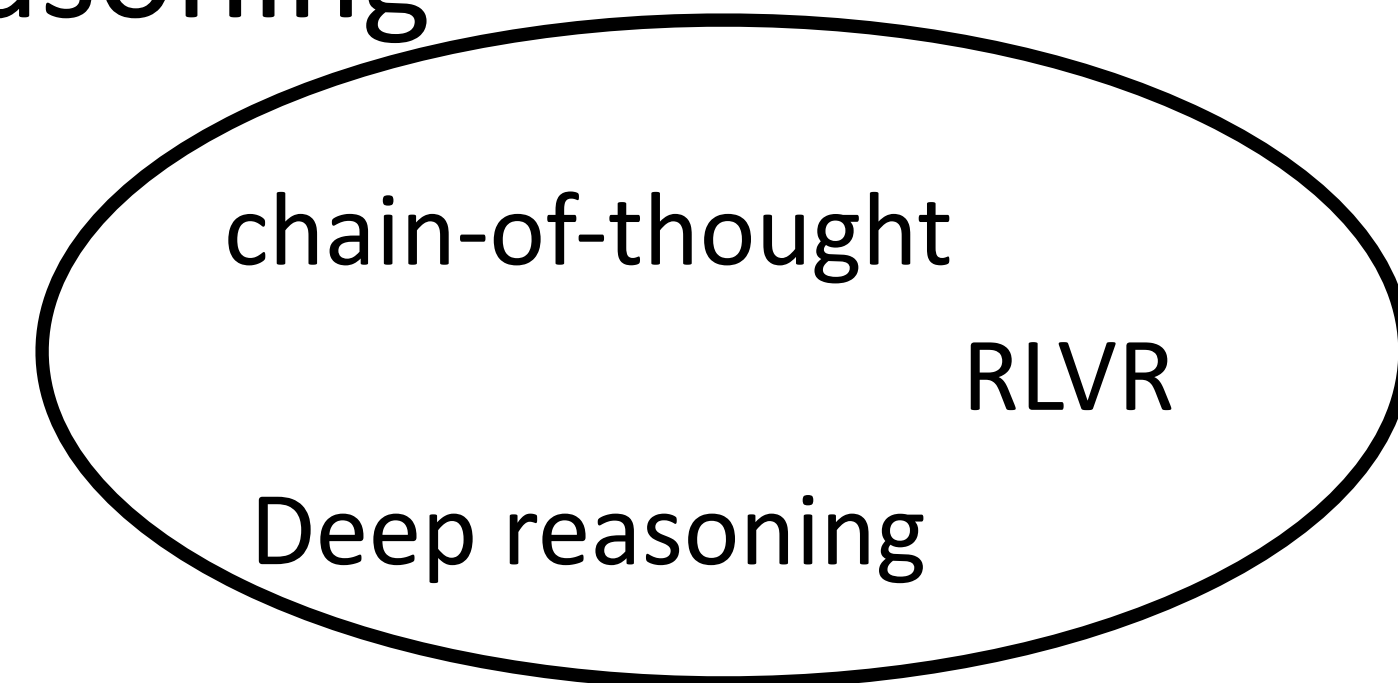
Assignment 3, 4

Agent



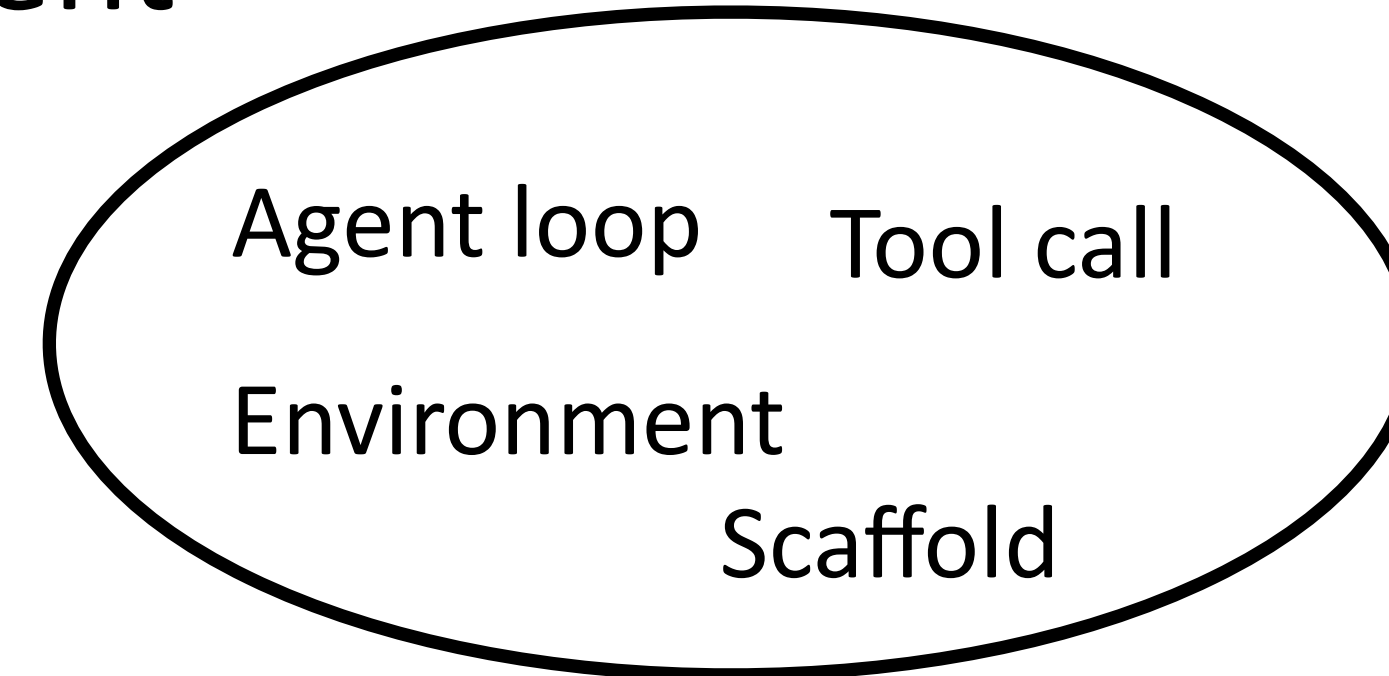
Review — Applications

Reasoning



Assignment 3, 4

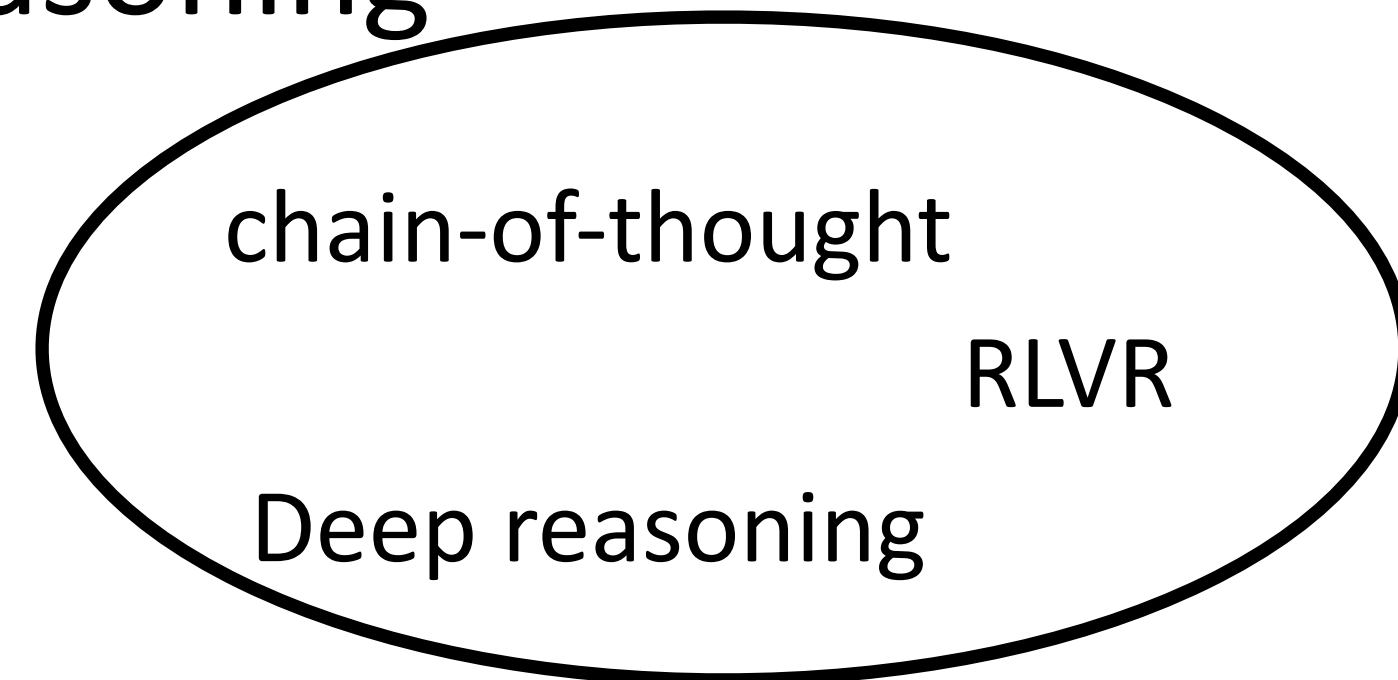
Agent



Group project + Assignment 4

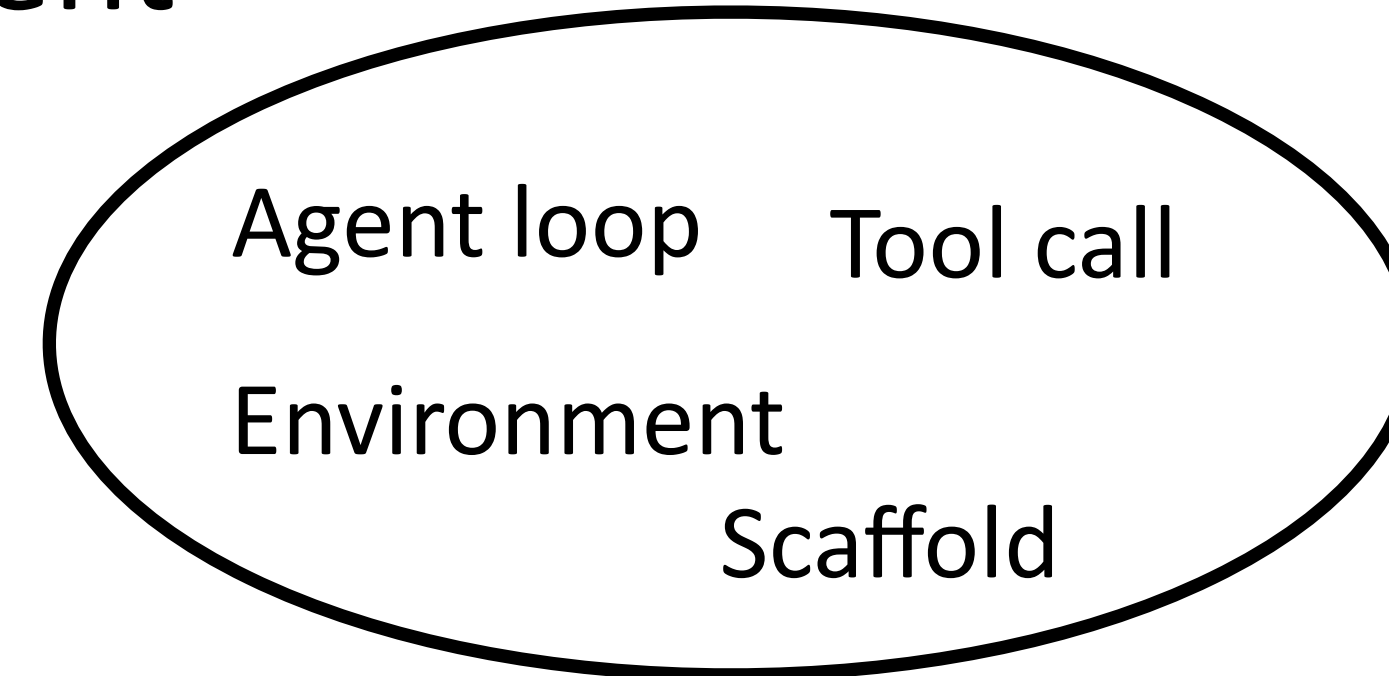
Review — Applications

Reasoning



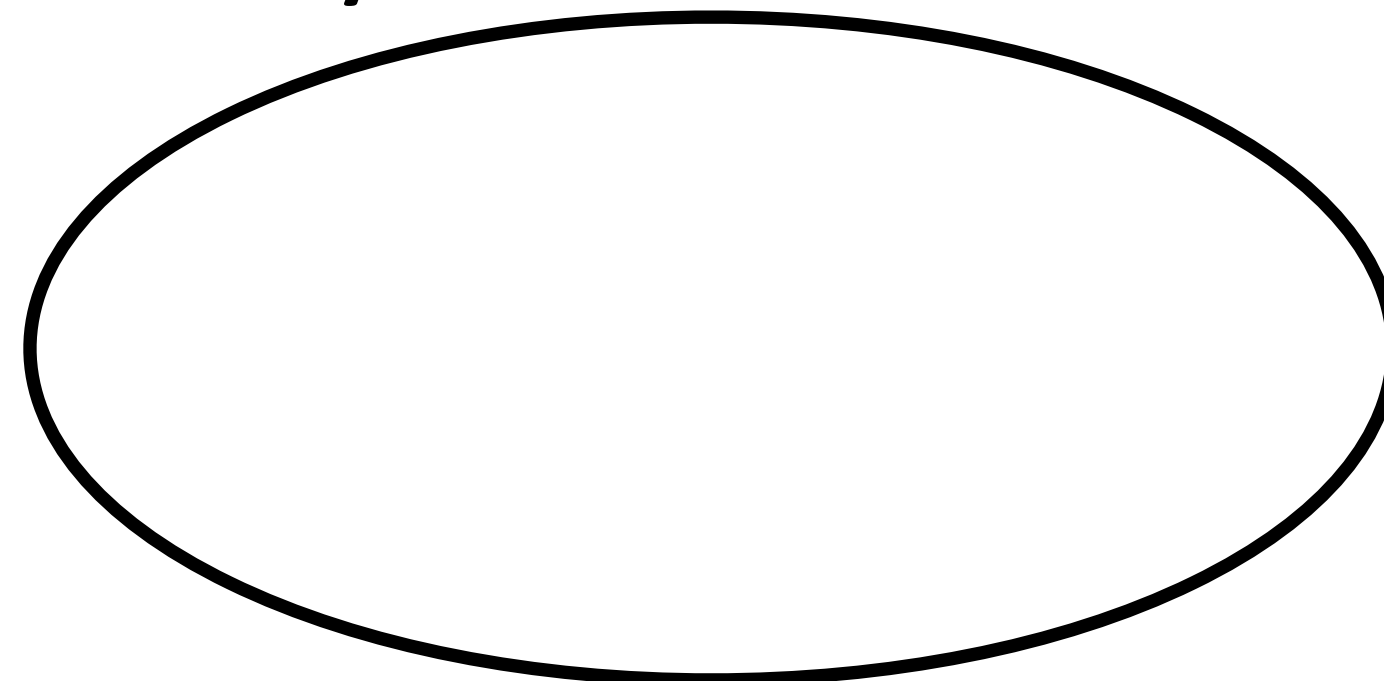
Assignment 3, 4

Agent



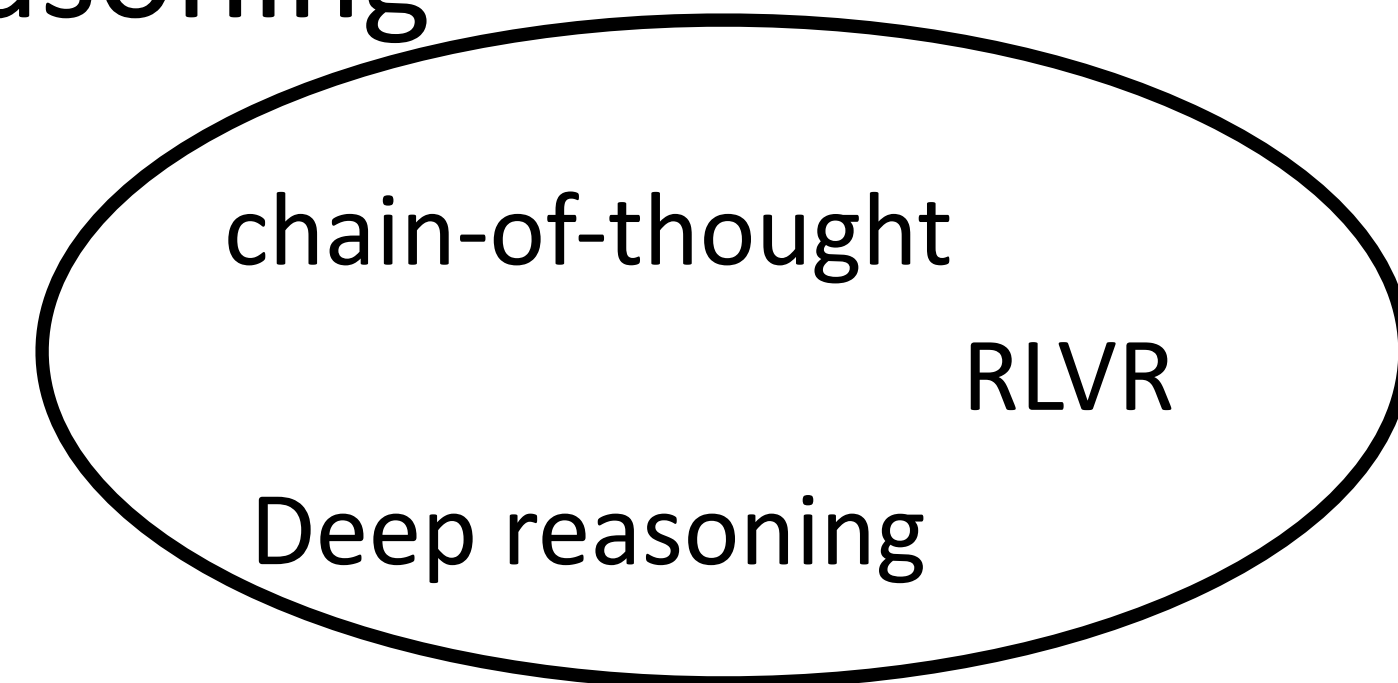
Group project + Assignment 4

Ethics/Safety



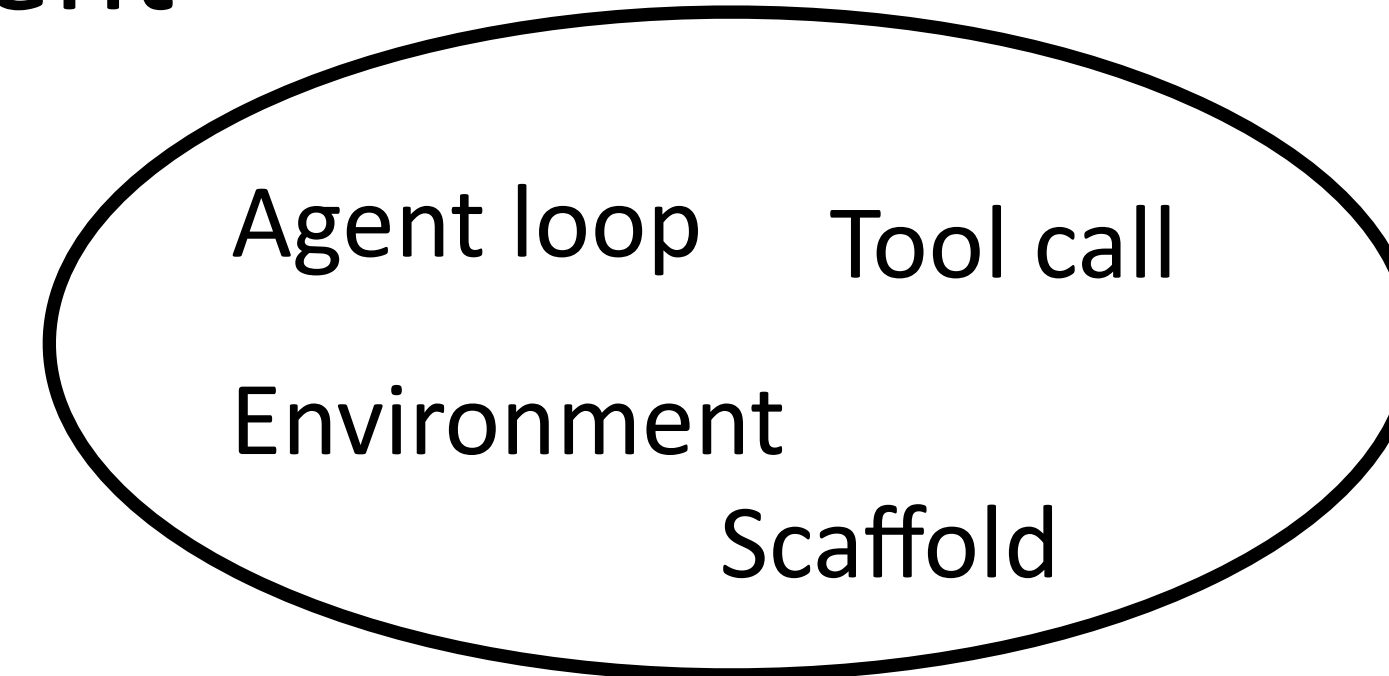
Review — Applications

Reasoning



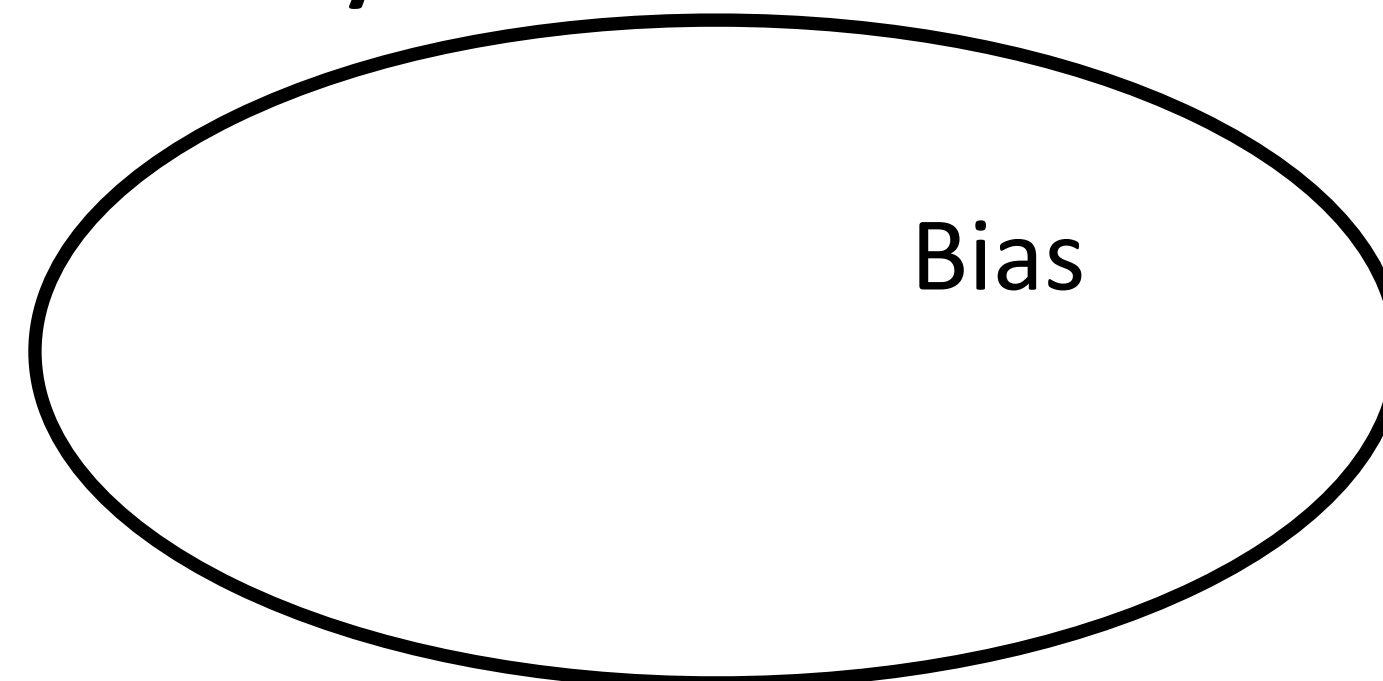
Assignment 3, 4

Agent



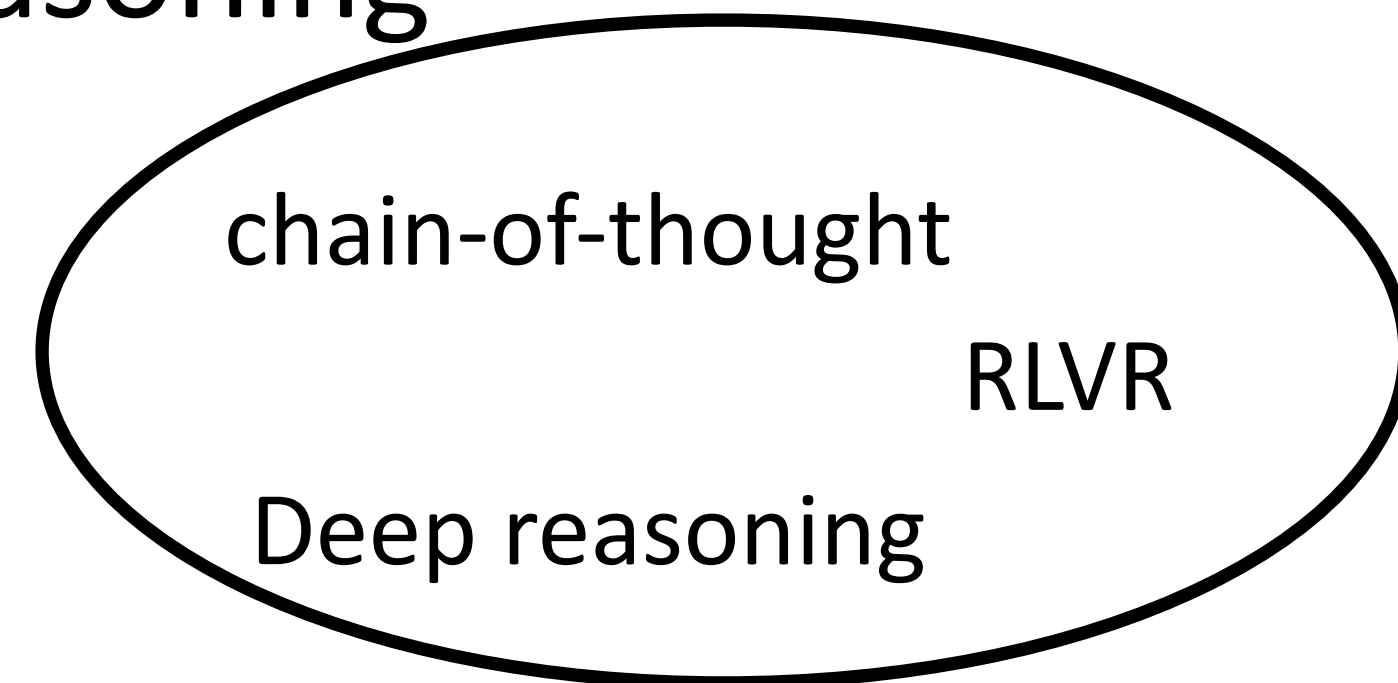
Group project + Assignment 4

Ethics/Safety



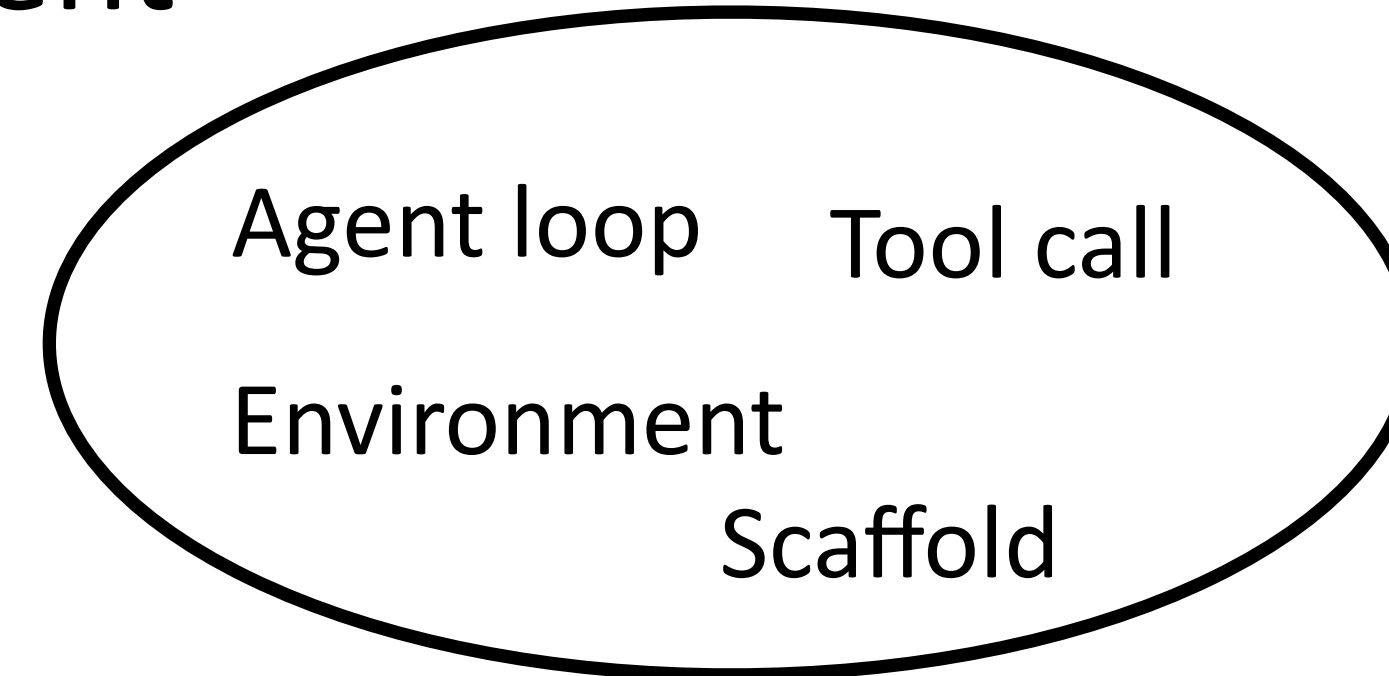
Review — Applications

Reasoning



Assignment 3, 4

Agent



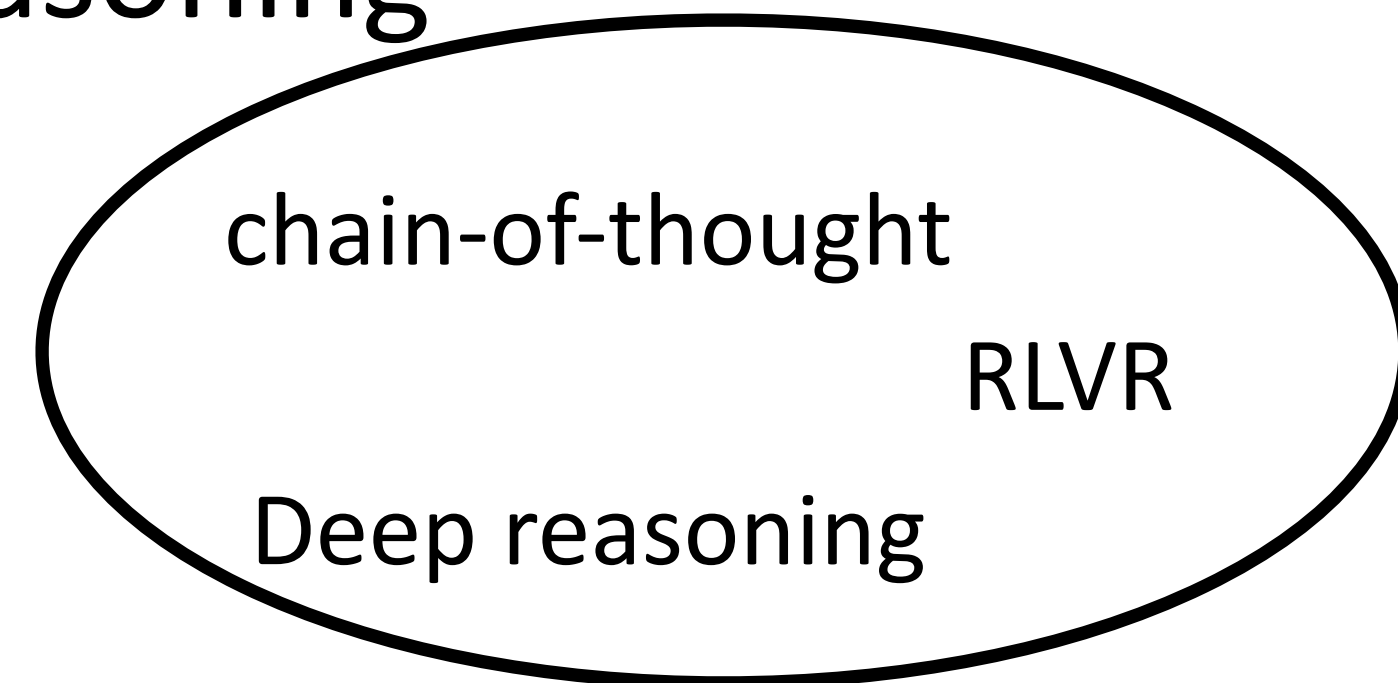
Group project + Assignment 4

Ethics/Safety



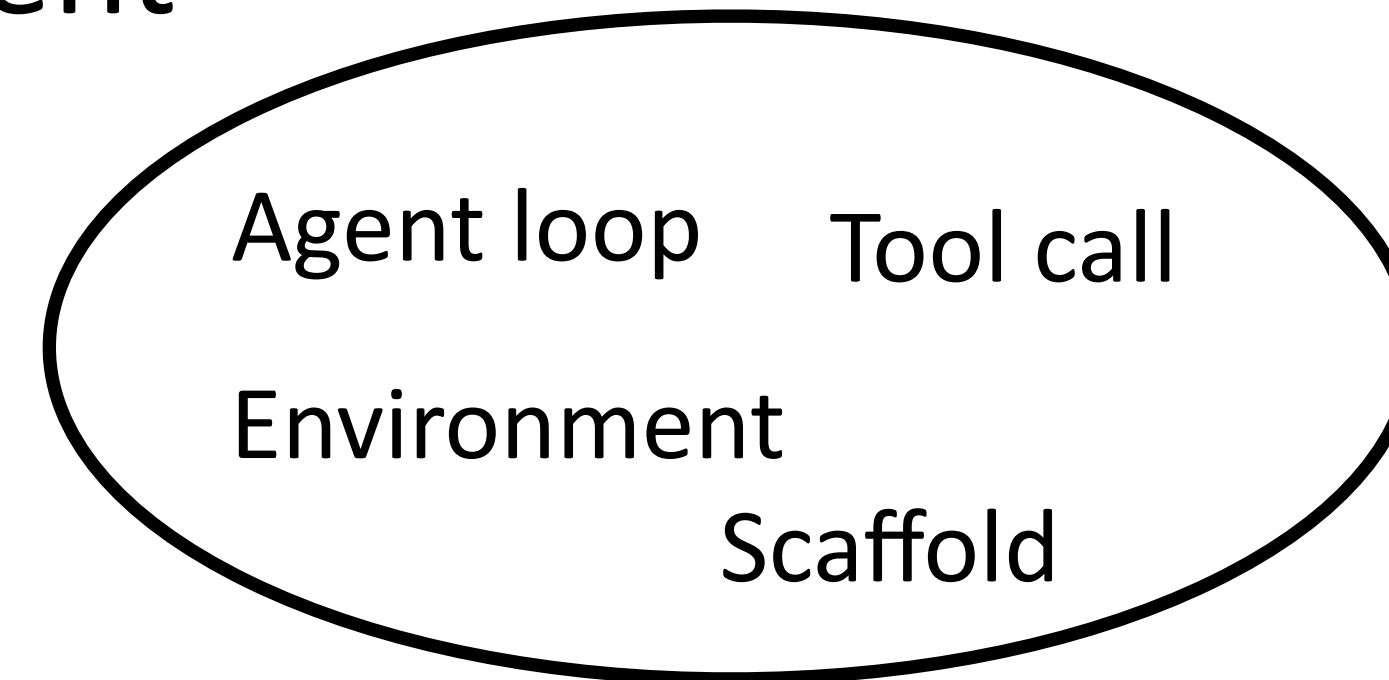
Review — Applications

Reasoning



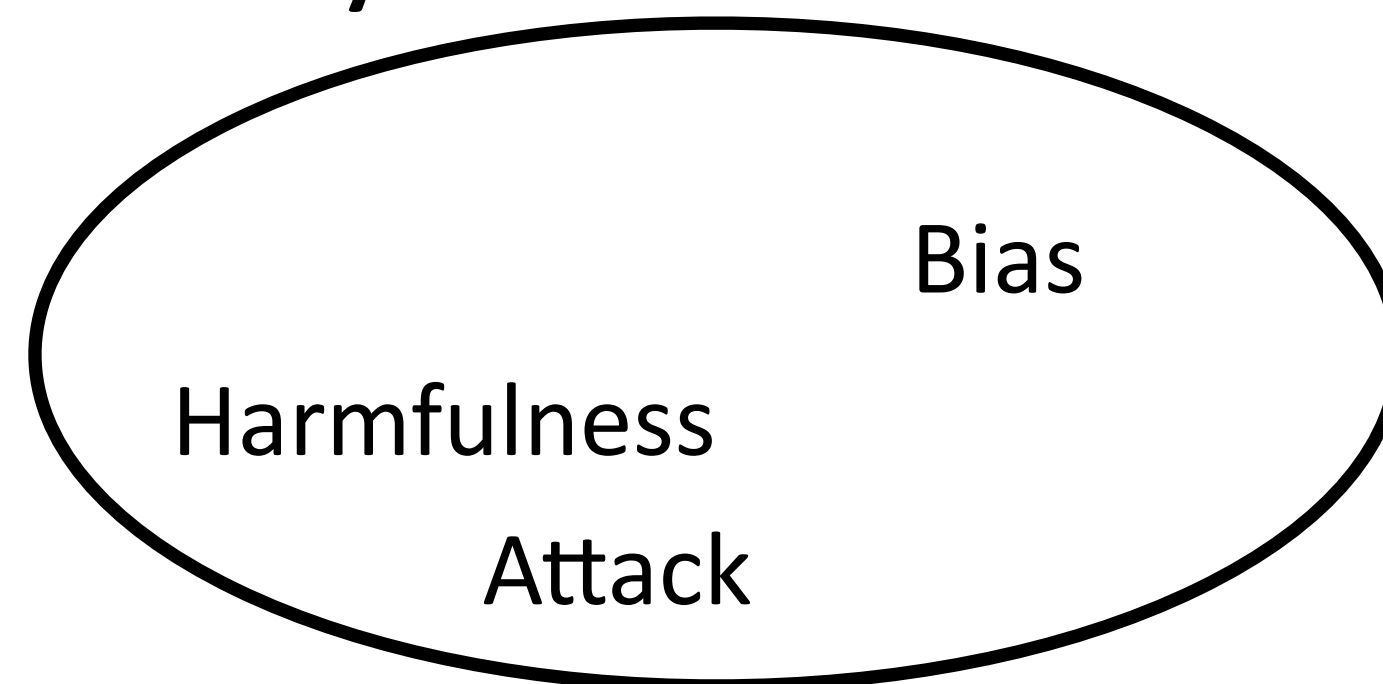
Assignment 3, 4

Agent



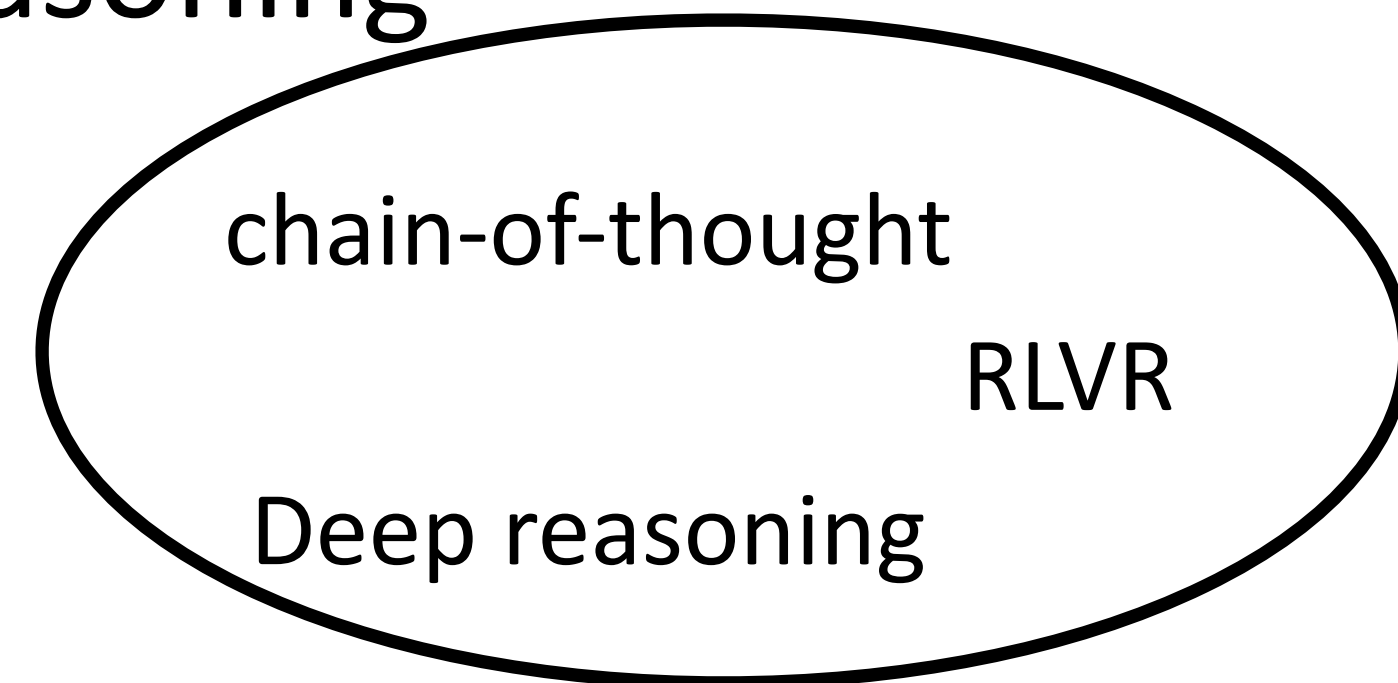
Group project + Assignment 4

Ethics/Safety



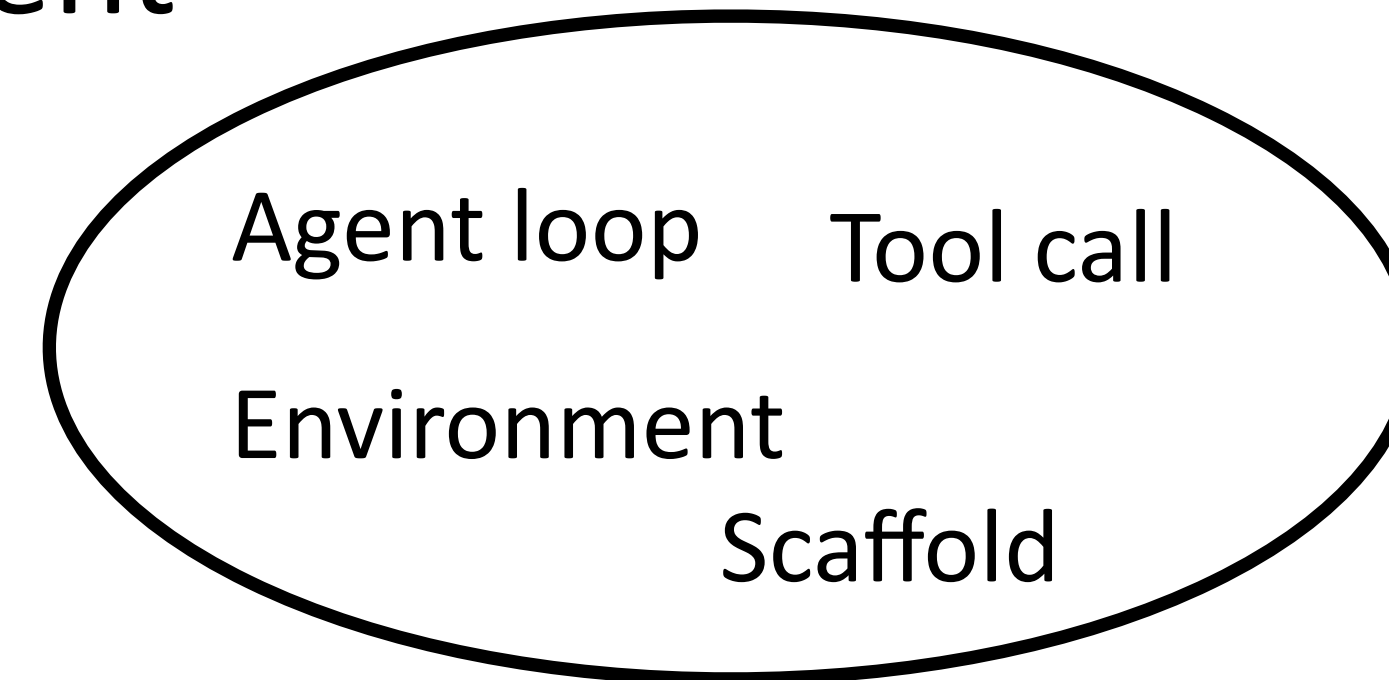
Review — Applications

Reasoning



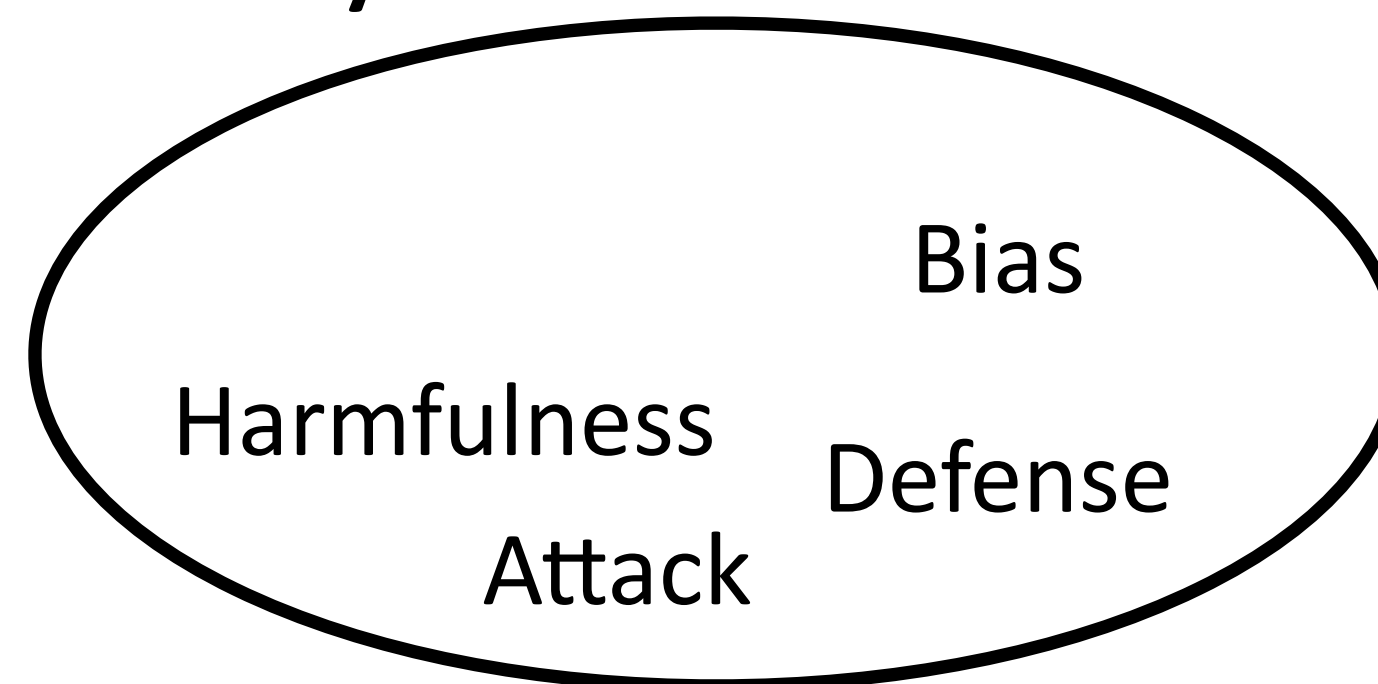
Assignment 3, 4

Agent



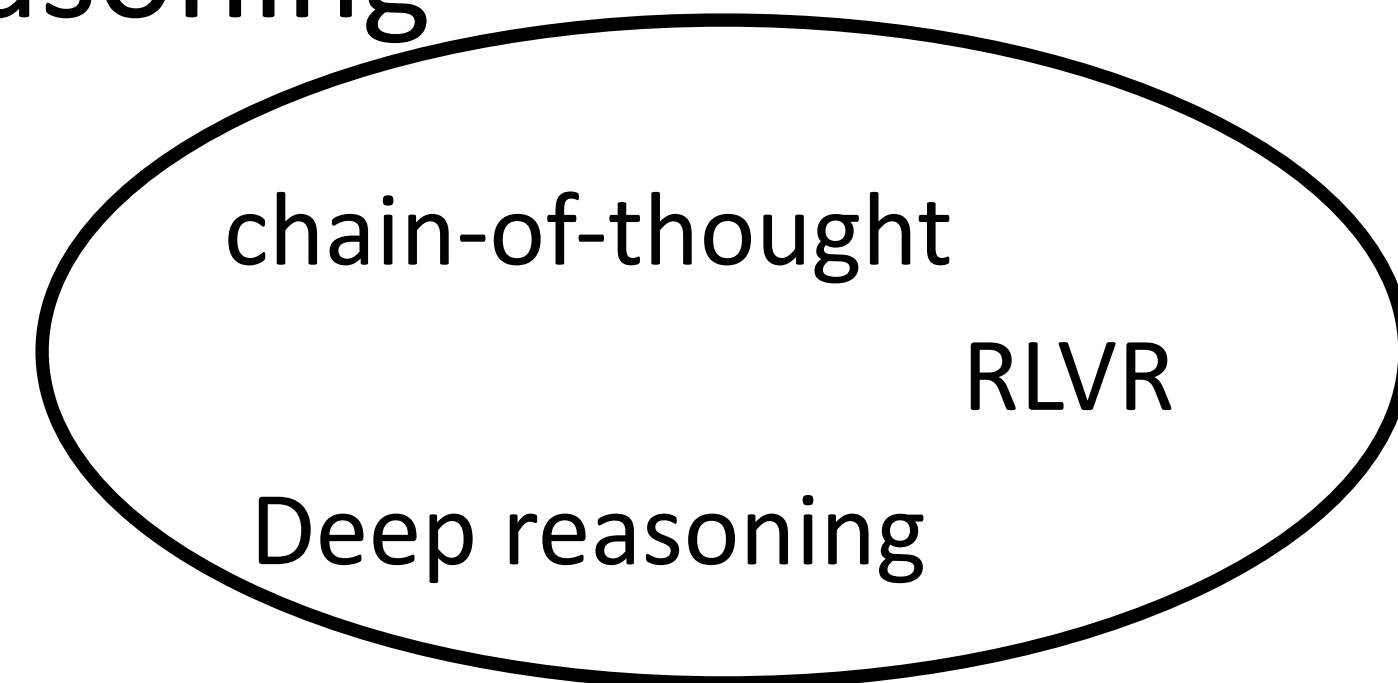
Group project + Assignment 4

Ethics/Safety



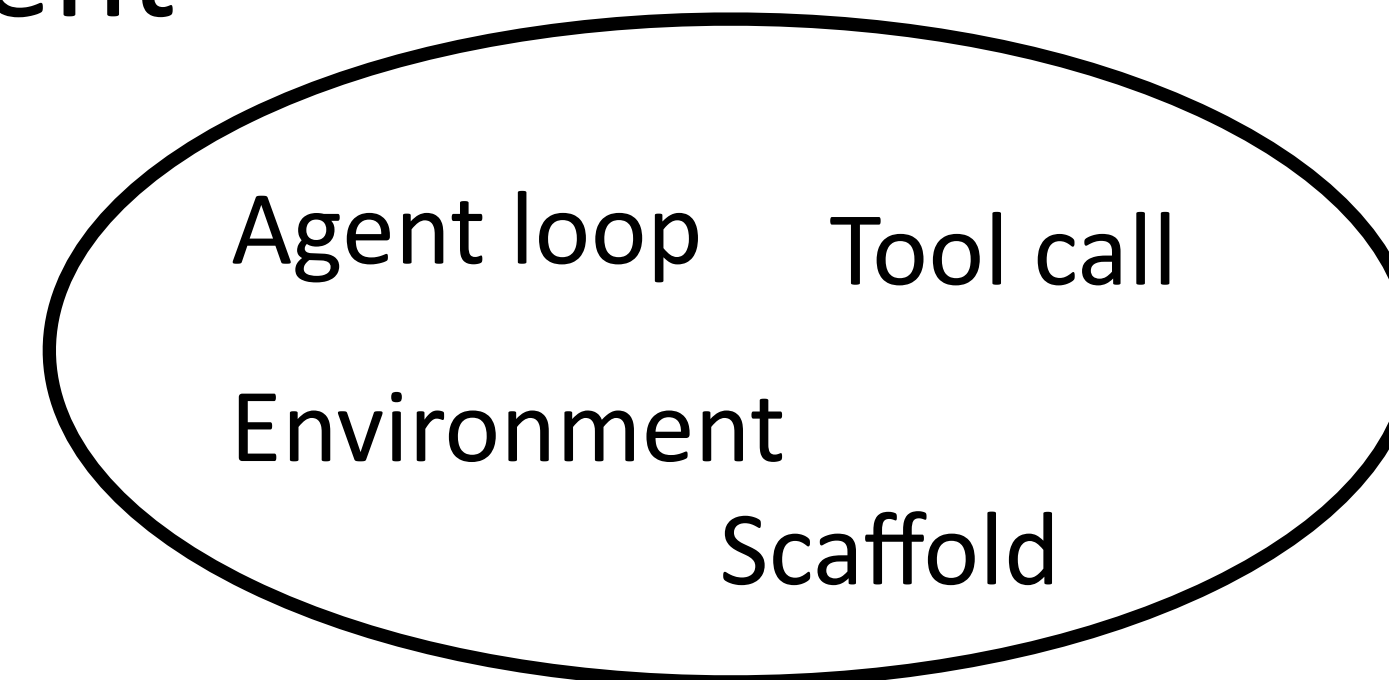
Review — Applications

Reasoning



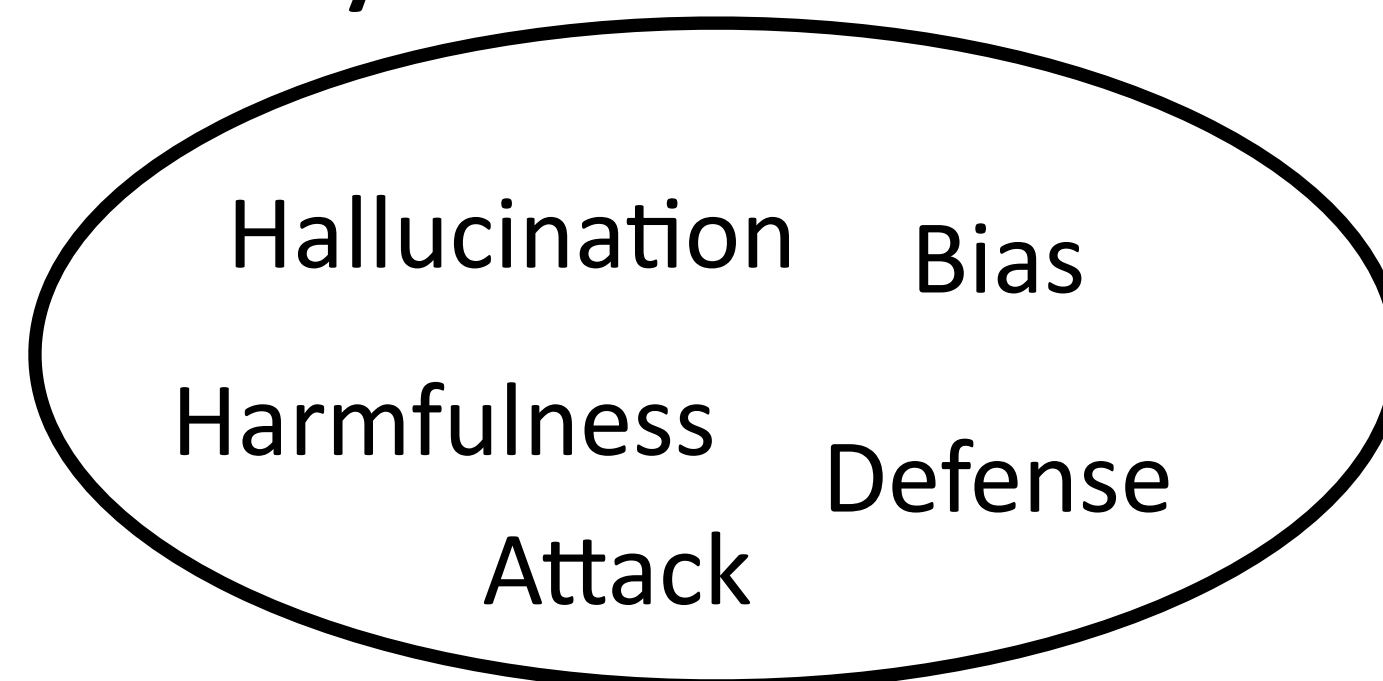
Assignment 3, 4

Agent



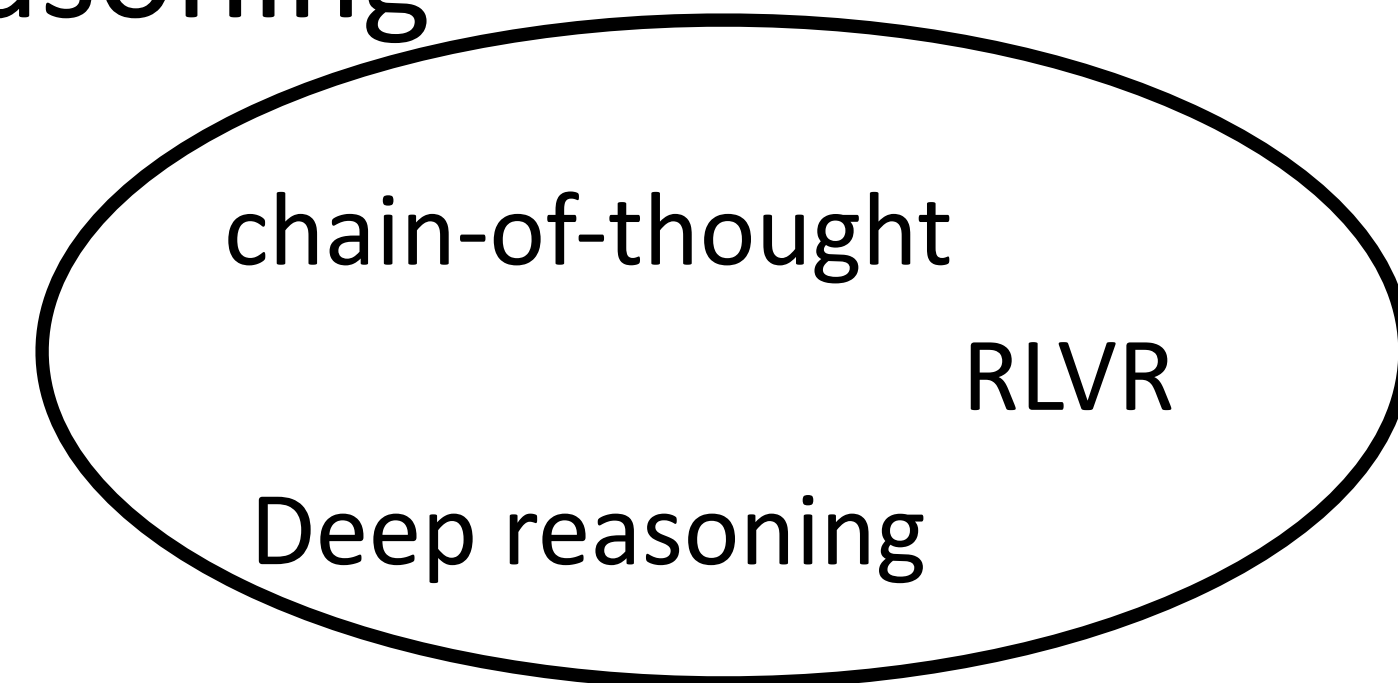
Group project + Assignment 4

Ethics/Safety



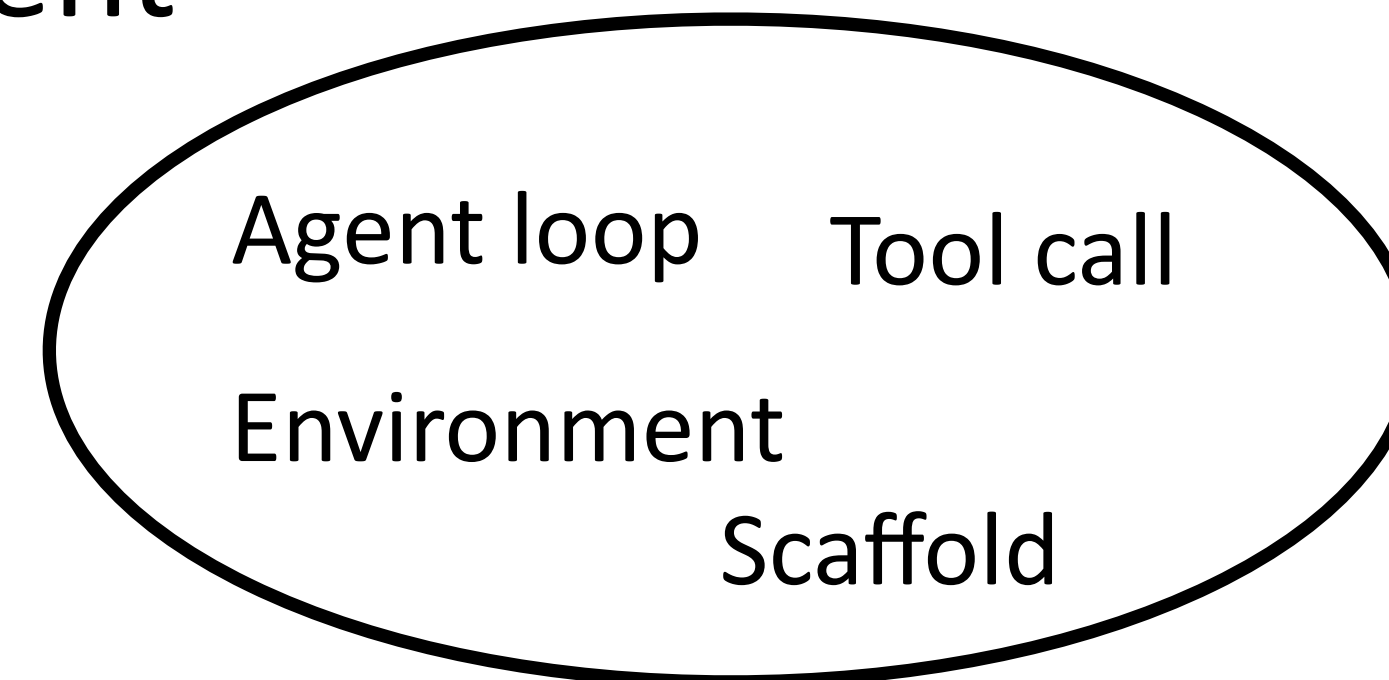
Review — Applications

Reasoning



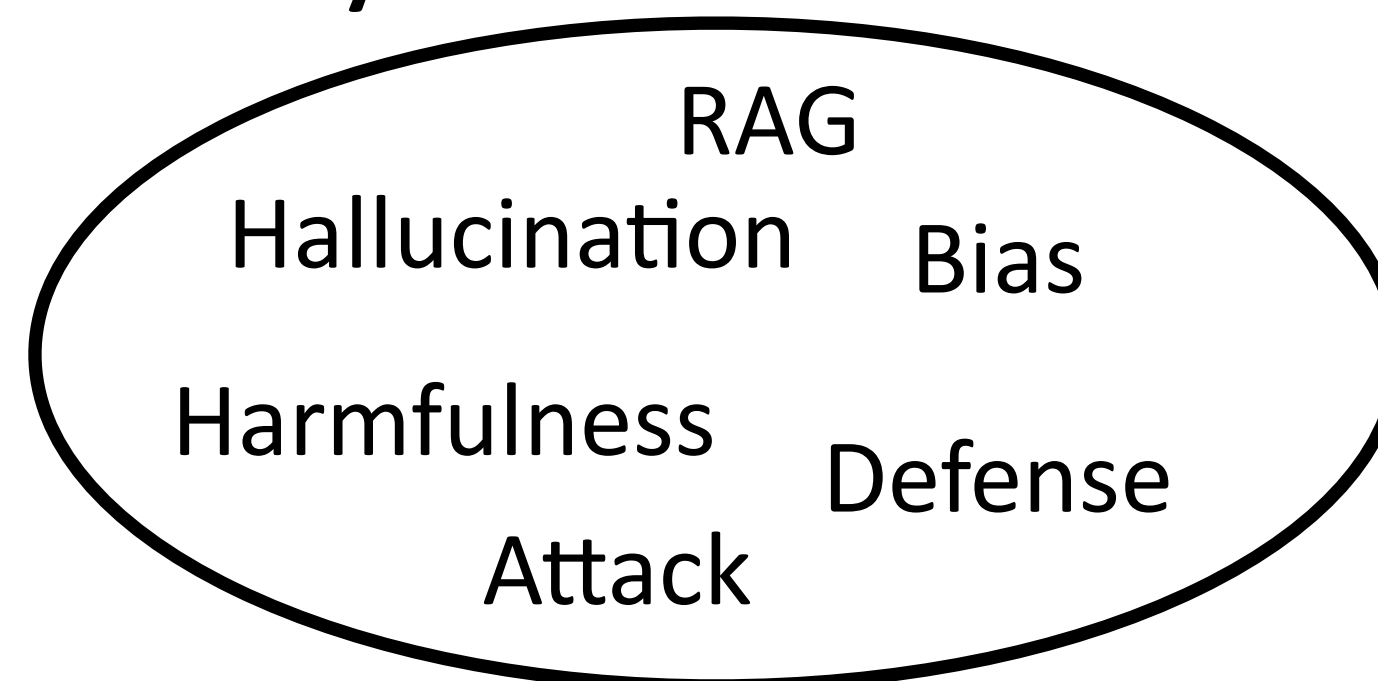
Assignment 3, 4

Agent



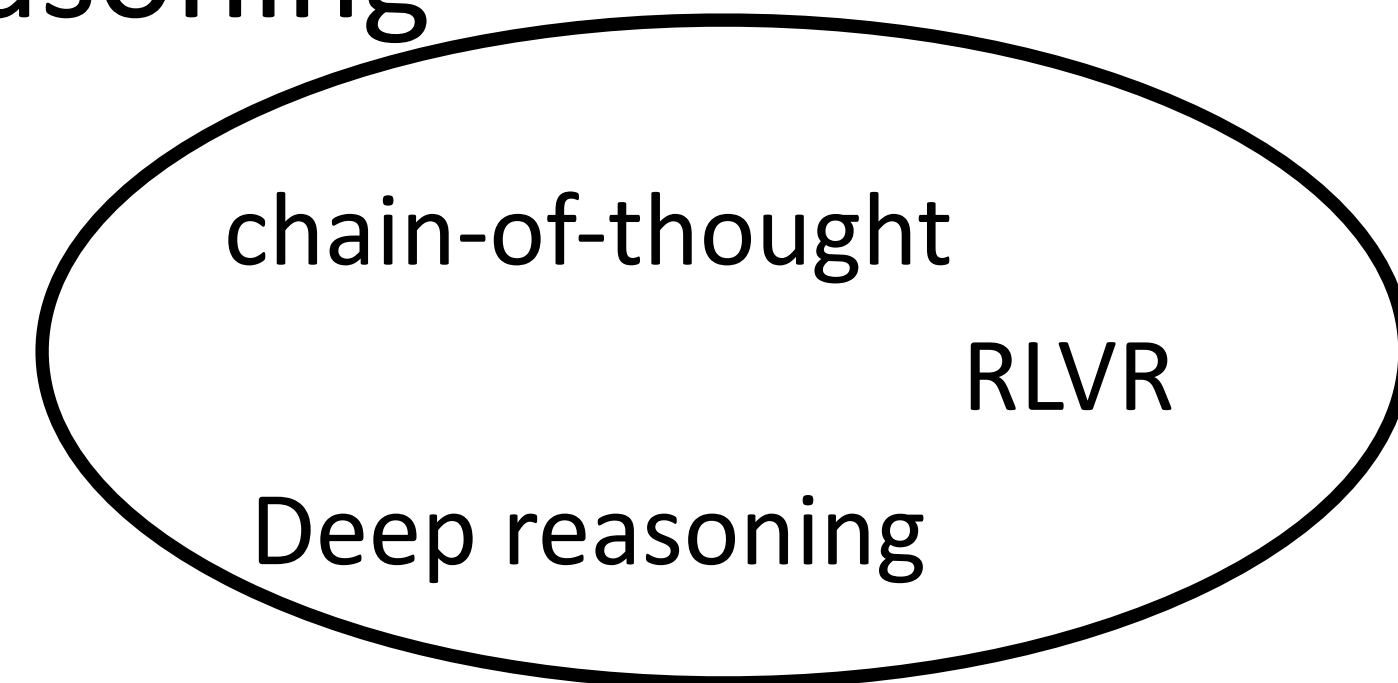
Group project + Assignment 4

Ethics/Safety



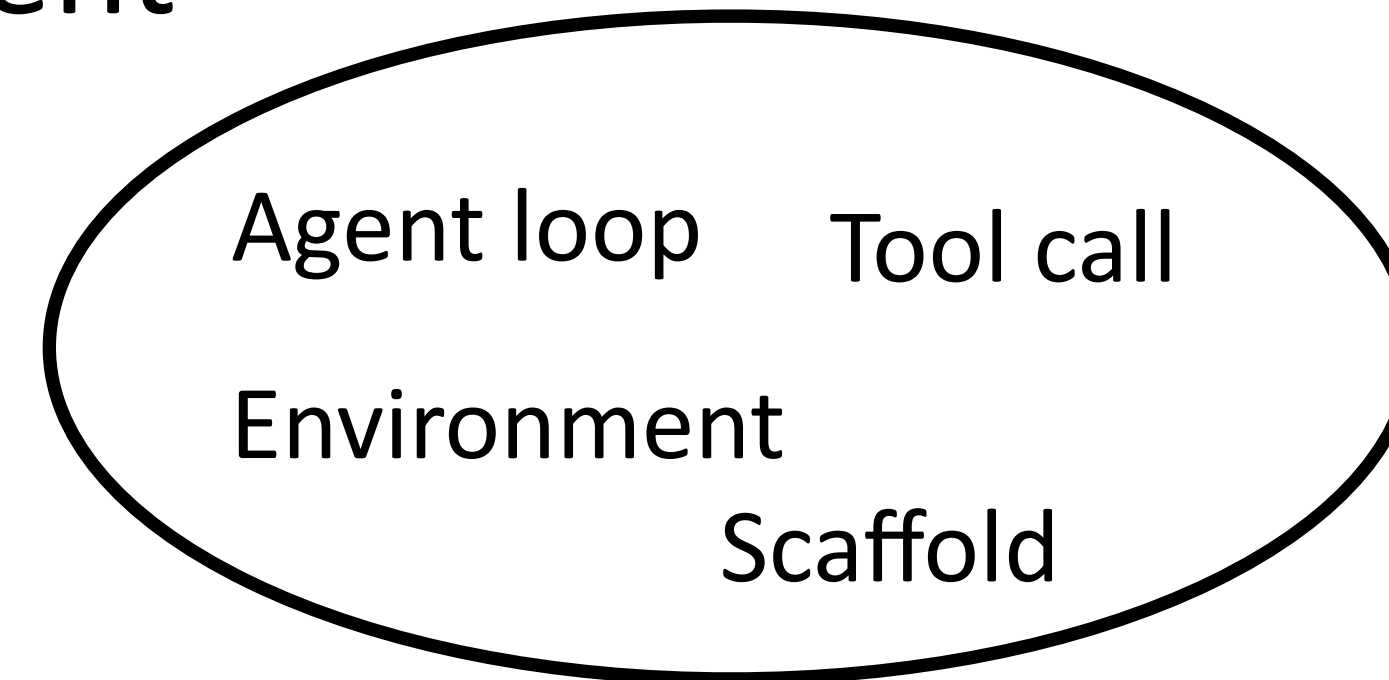
Review — Applications

Reasoning



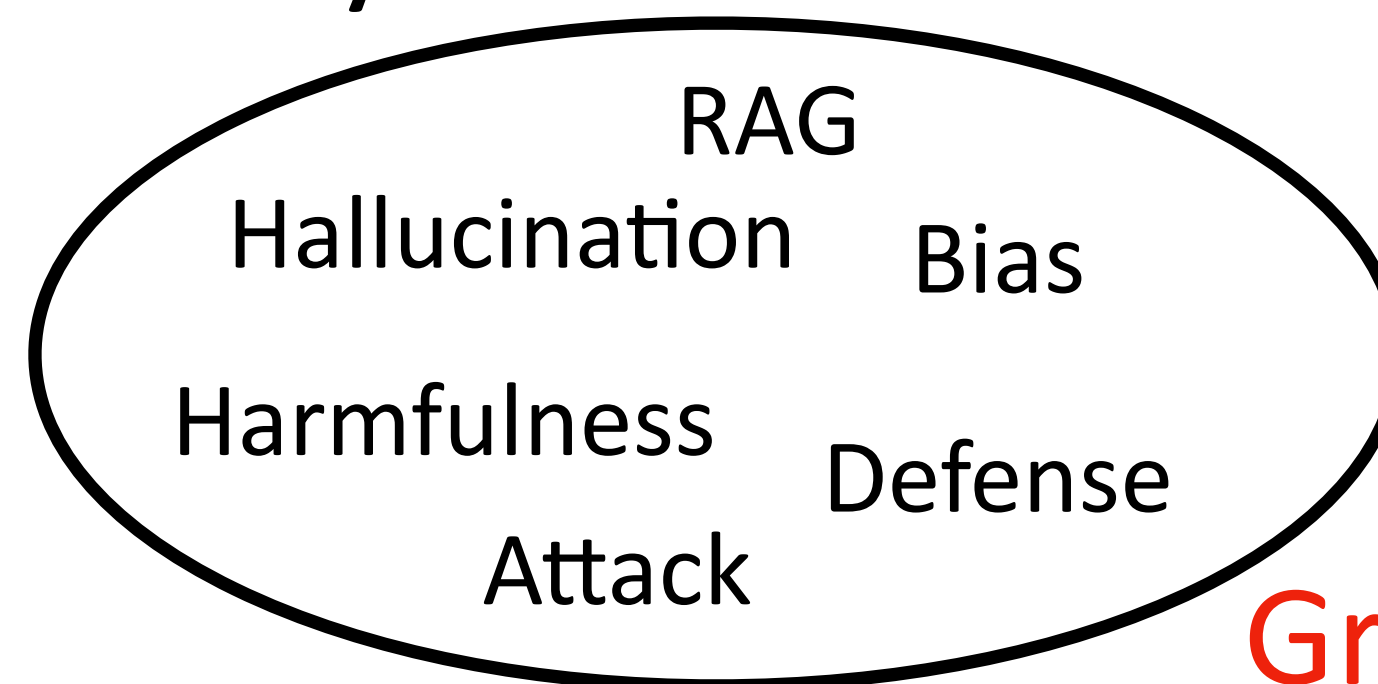
Assignment 3, 4

Agent



Group project + Assignment 4

Ethics/Safety



Group project

About Large Language Models

About Large Language Models

1. Learning from practice (e.g., programming assignments, projects) is the most important

About Large Language Models

1. Learning from practice (e.g., programming assignments, projects) is the most important
2. This field is evolving very fast (like many things change every 6 months), some of the course contents may get out-of-date, need to catch up things quickly

About Large Language Models

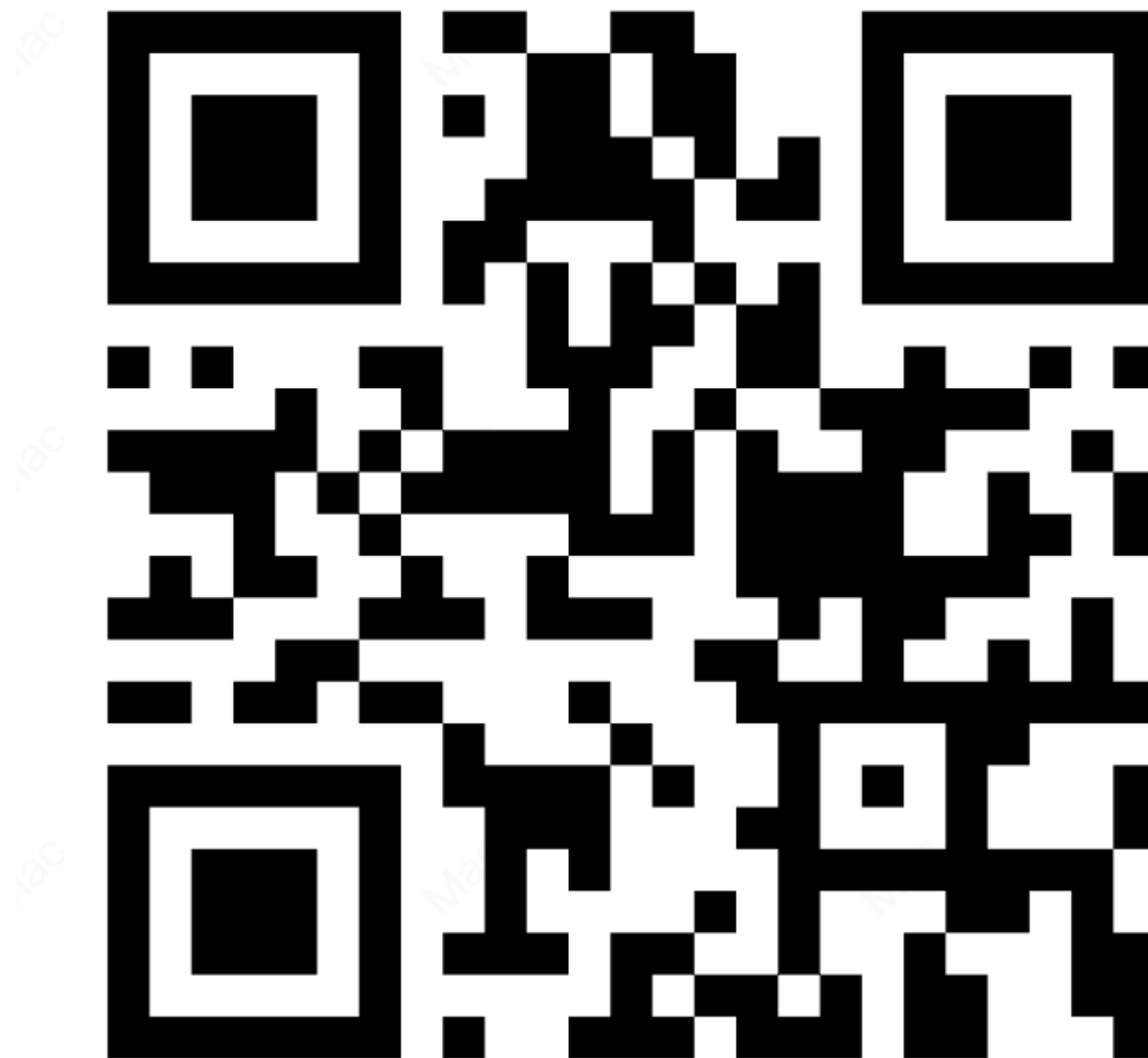
1. Learning from practice (e.g., programming assignments, projects) is the most important
2. This field is evolving very fast (like many things change every 6 months), some of the course contents may get out-of-date, need to catch up things quickly
3. Always embrace new technologies / applications

About Large Language Models

1. Learning from practice (e.g., programming assignments, projects) is the most important
2. This field is evolving very fast (like many things change every 6 months), some of the course contents may get out-of-date, need to catch up things quickly
3. Always embrace new technologies / applications
4. LLMs allow you to build amazing stuff quickly that can make actual impacts

Student Feedback Questionnaire (SFQ)

Anonymous and not mandatory, ddl tomorrow (Nov 29)



<https://sfq-survey.ust.hk>

Thank You for taking the course!