

Multimodal AI

Lecture 14.1 – Advanced Topics

Paul Liang

Assistant Professor

MIT Media Lab & MIT EECS



<https://pliang279.github.io>

ppliang@mit.edu

 [@pliang279](https://twitter.com/pliang279)



Assignments for This Coming Week

Please fill out course evaluations and give us feedback!

HW5 due this Friday May 8.

For project:

- Make sure to meet with myself and TAs this week if you need advice.
- Presentations next Tuesday May 12.
- Final report due Tuesday May 19.

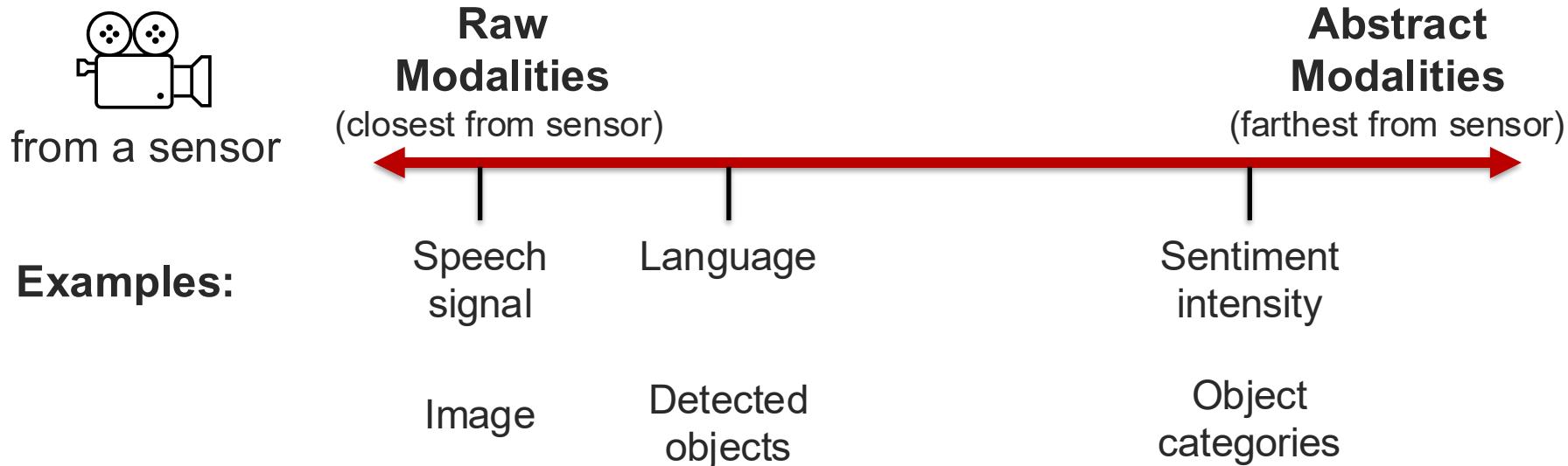
This Week: Recap and Advanced Topics

1. Recap of course material
2. Native multimodal models and mixture of experts
3. Advanced fusion techniques
4. Self-evolving multimodal agents
5. Extending human senses – touch, smell, and taste
6. Human-AI interaction and safety

What is a Modality?

Modality

Modality refers to the way in which something expressed or perceived.



What is Multimodal?

A dictionary definition...

Multimodal: with multiple modalities

A research-oriented definition...

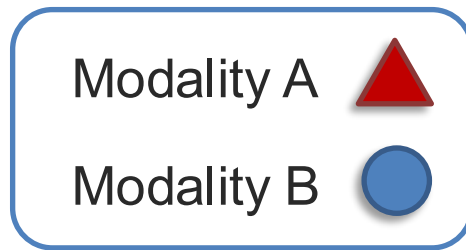
***Multimodal* is the science of**

heterogeneous and interconnected data

Connected + Interacting

Heterogeneous Modalities

Information in different modalities shows diverse qualities, structures, & representations.



Homogeneous Modalities
(with similar qualities)

Heterogeneous Modalities
(with diverse qualities)



Images
from 2
cameras



Text from
2 different
languages



Language
and vision



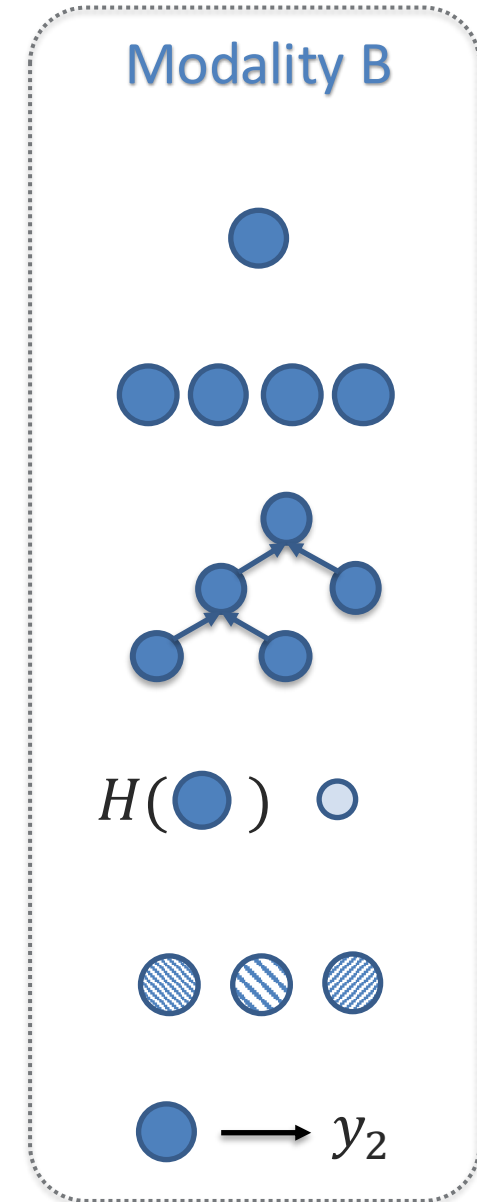
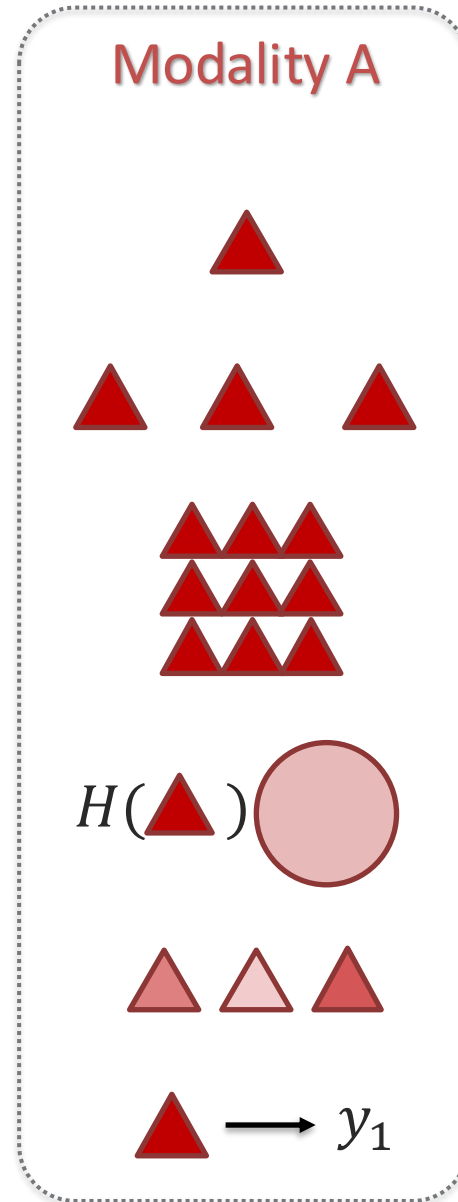
Language
and sensors

Examples:

Abstract modalities are more likely to be homogeneous

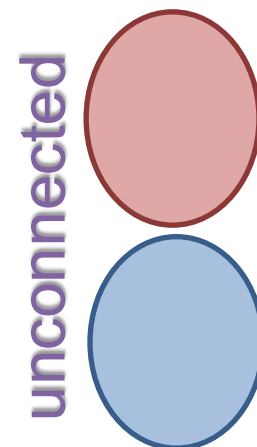
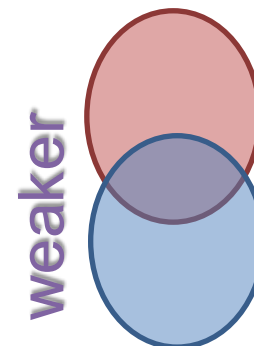
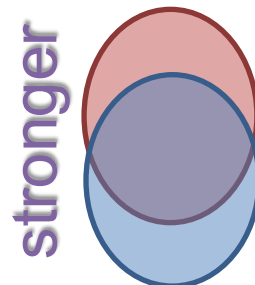
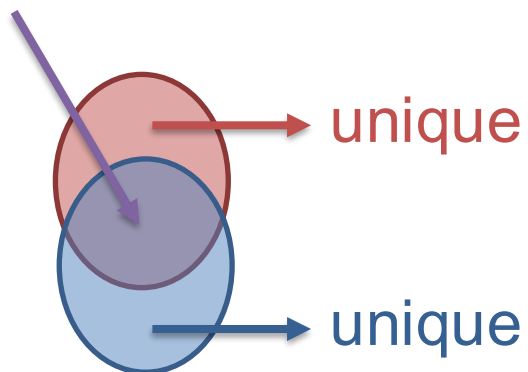
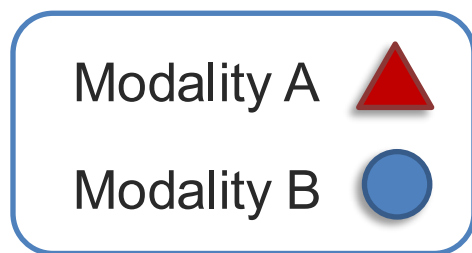
Modality Profile

- 1 **Element representations:**
Discrete, continuous, granularity
- 2 **Element distributions:**
Density, frequency
- 3 **Structure:**
Temporal, spatial, latent, explicit
- 4 **Information:**
Abstraction, entropy
- 5 **Noise:**
Uncertainty, noise, missing data
- 6 **Relevance:**
Task, context dependence



Connected Modalities

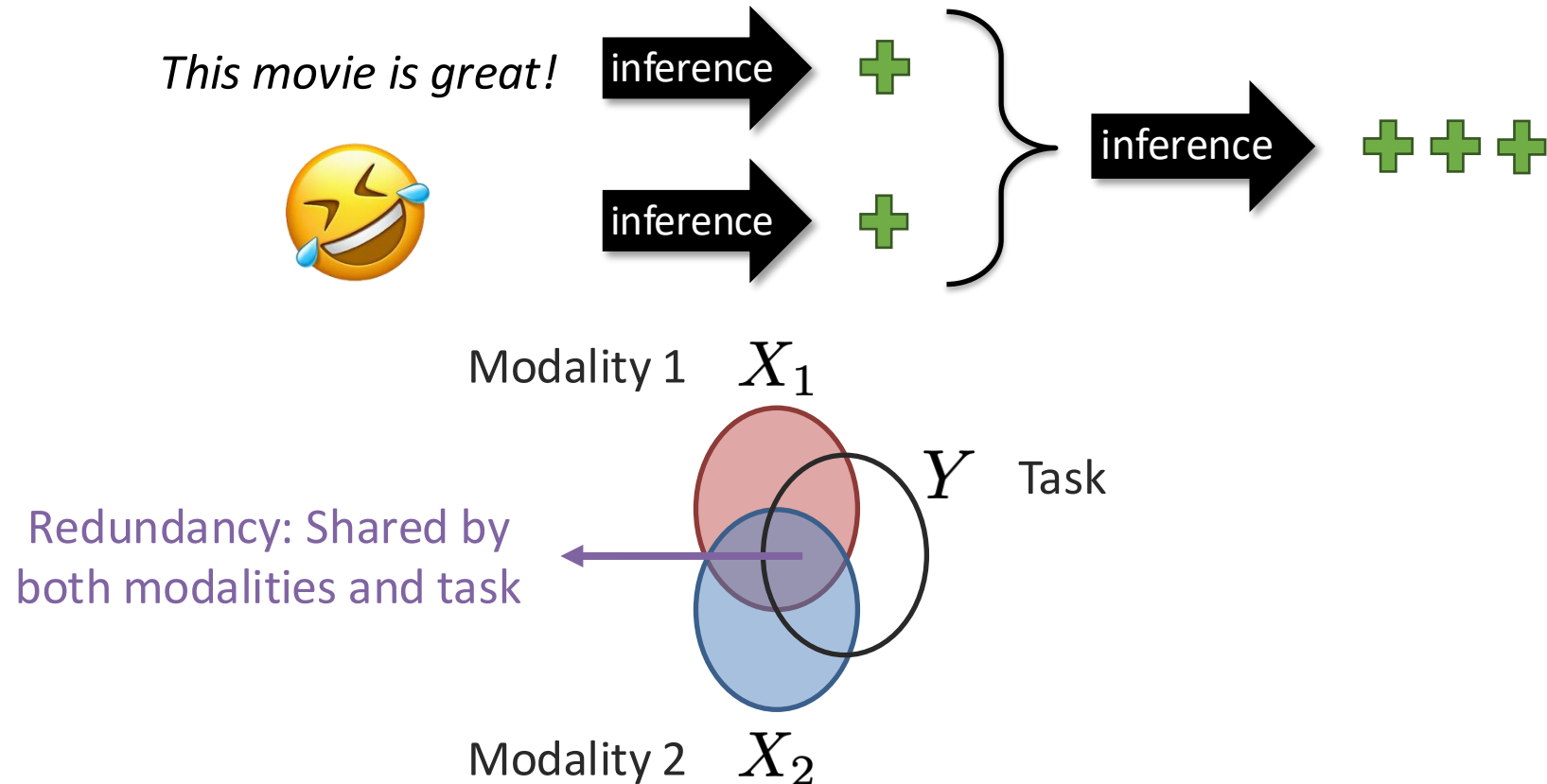
Shared information that relates modalities



*A teacup on the right of a laptop
in a clean room.*

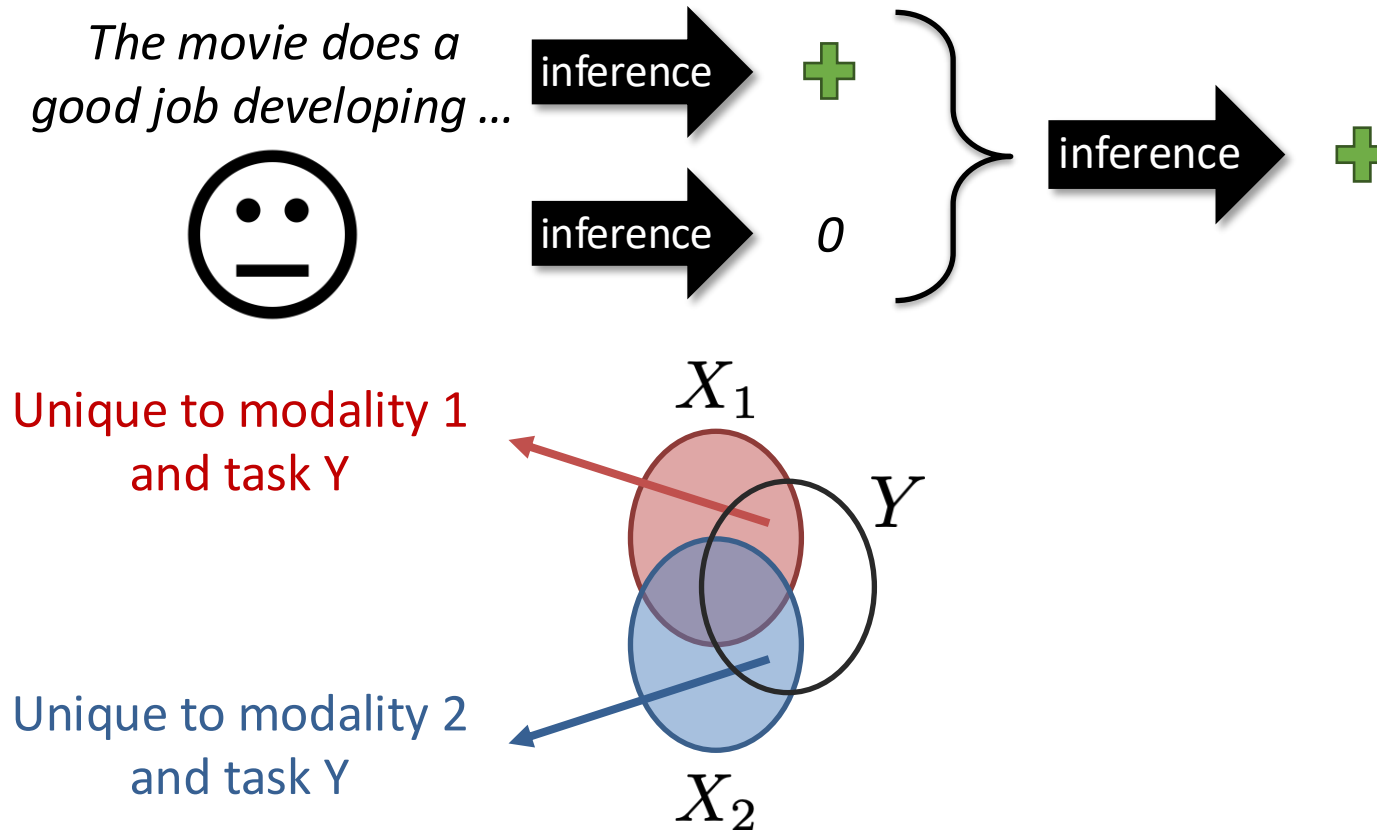
Interacting Modalities

Interactions: How modalities *combine* to provide information for a task.



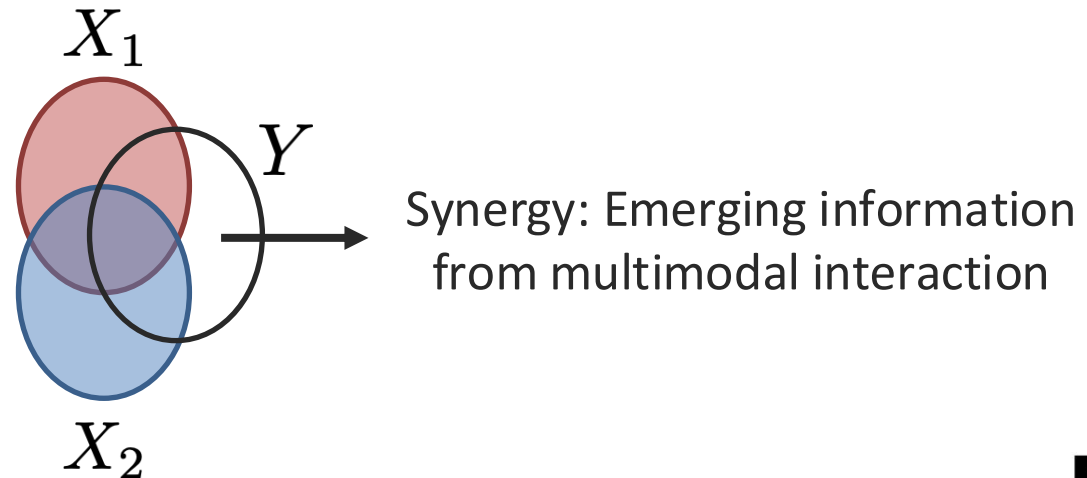
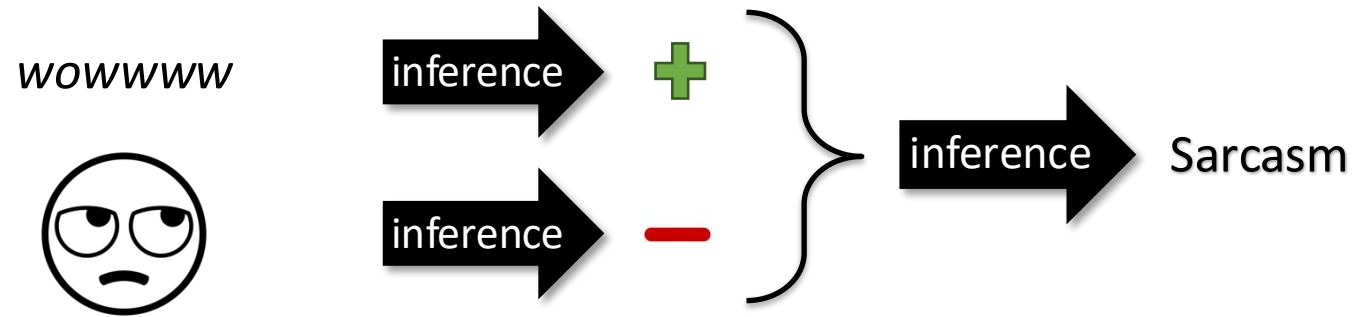
Interacting Modalities

Interactions: How modalities *combine* to provide information for a task.



Interacting Modalities

Interactions: How modalities *combine* to provide information for a task.



*What is
Multimodal?*



Why is it hard?



What is next?

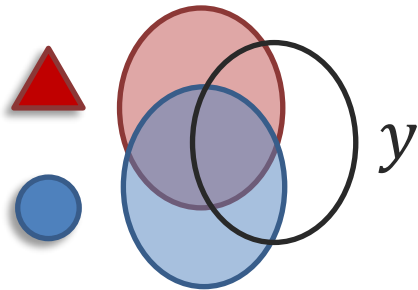
Heterogeneous



Connected

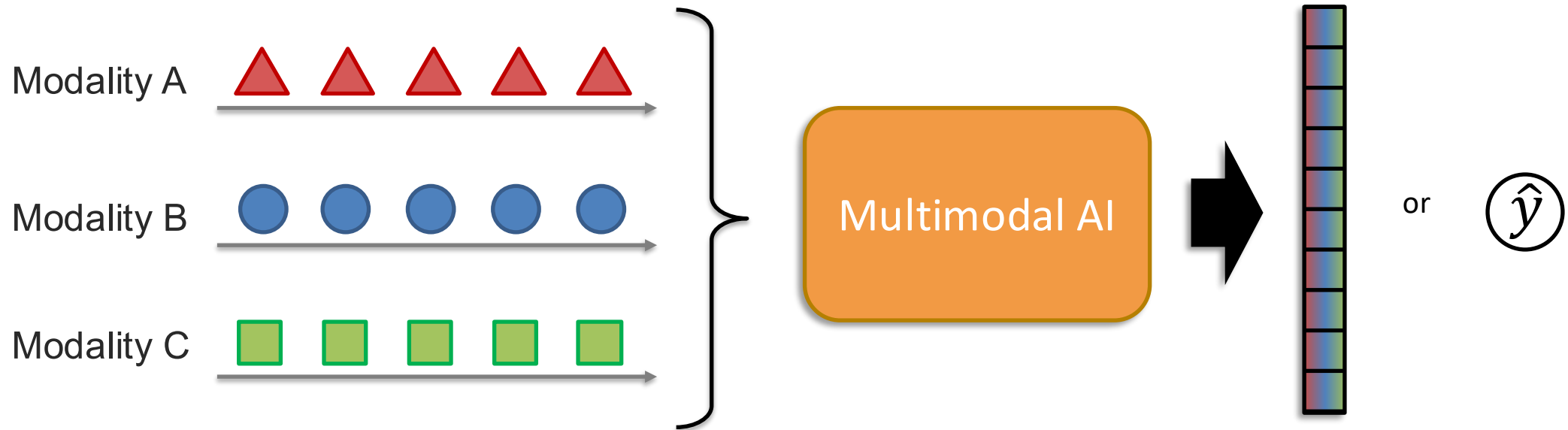


Interacting

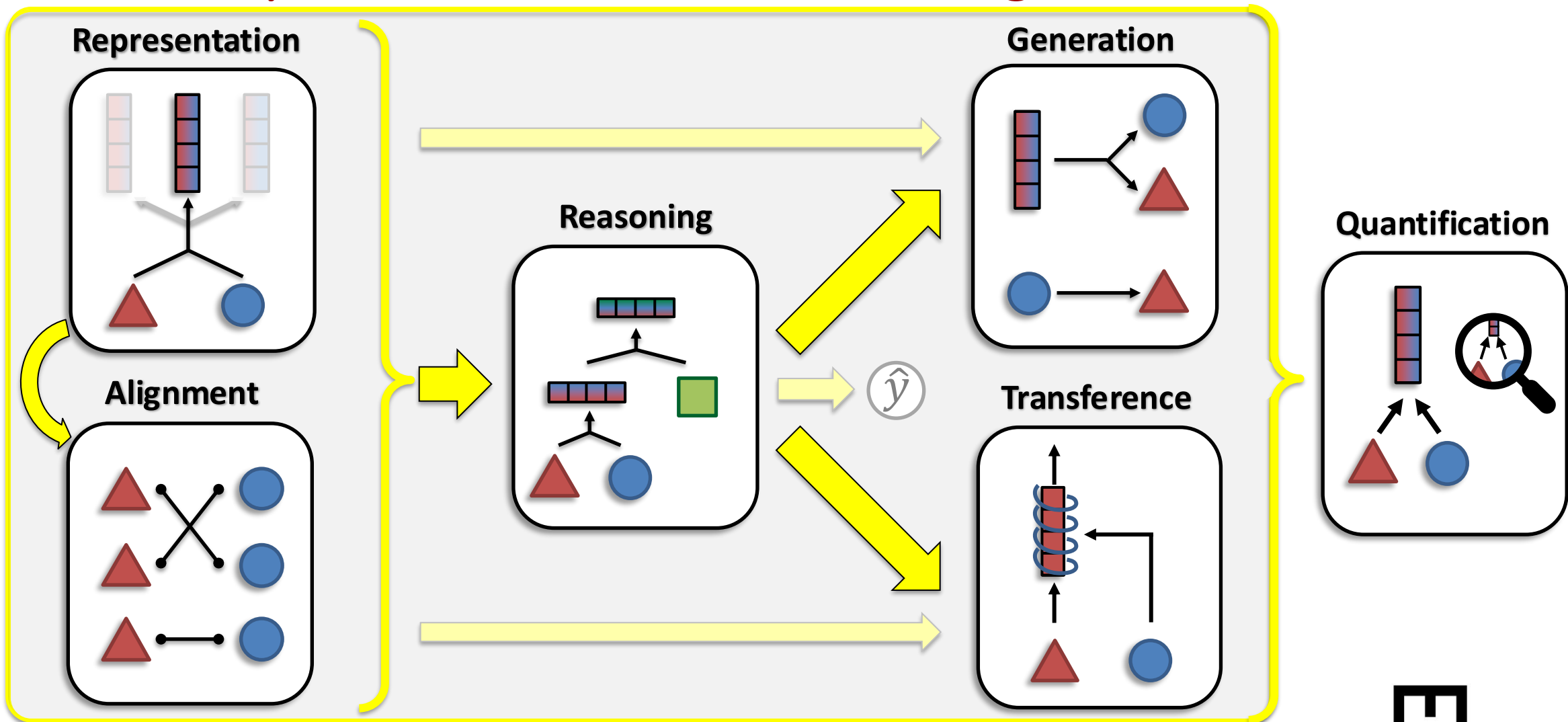


**Multimodal is the scientific
study of heterogeneous and
interconnected data 😊**

Multimodal AI Challenges



Summary of Core Multimodal Challenges

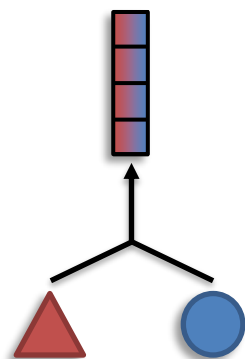


Challenge 1: Representation

Definition: Learning representations that reflect cross-modal interactions between individual elements, across different modalities.

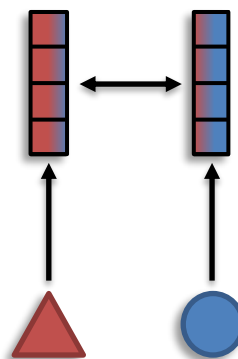
Sub-challenges:

Fusion



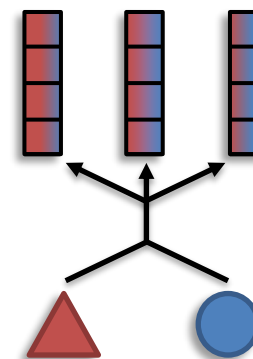
modalities \gt # representations

Coordination



modalities $=$ # representations

Fission



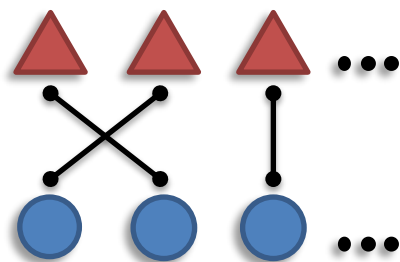
modalities \lt # representations

Challenge 2: Alignment

Definition: Identifying and modeling cross-modal connections between all elements of multiple modalities, building from the data structure.

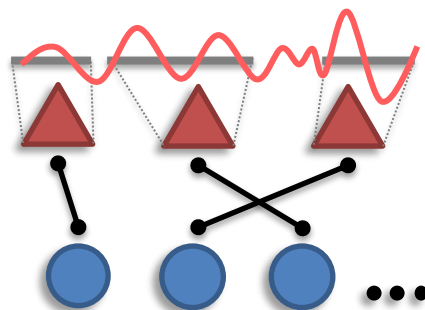
Sub-challenges:

Discrete connections



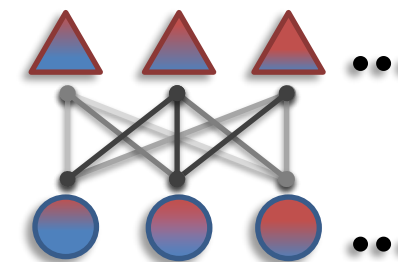
Explicit alignment
(e.g., grounding)

Continuous alignment



Granularity of
individual elements

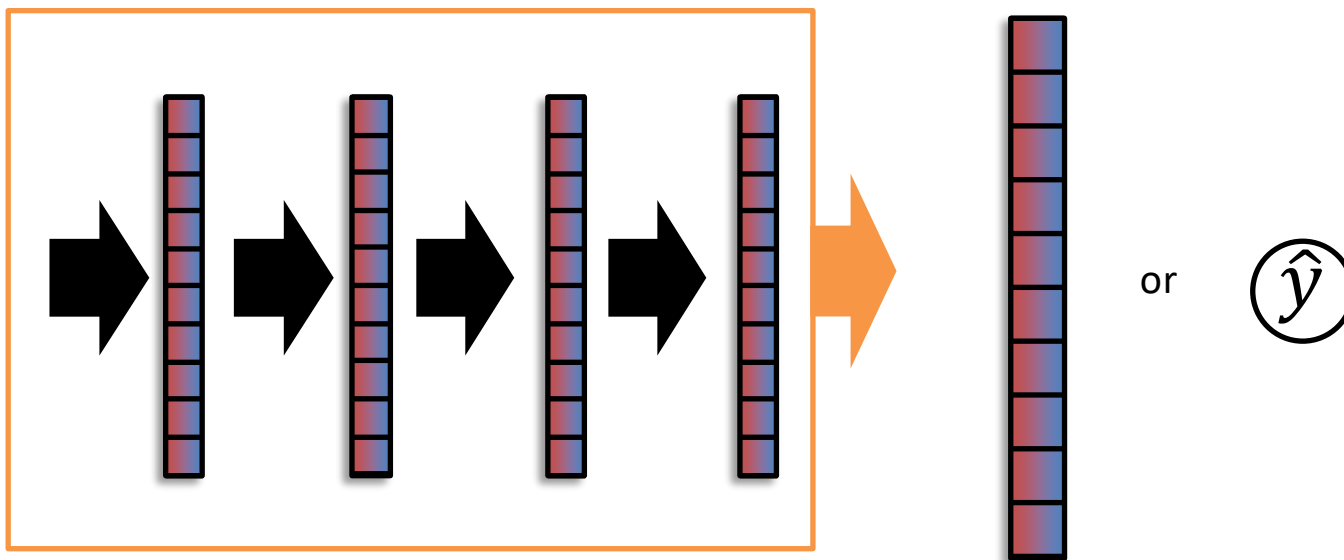
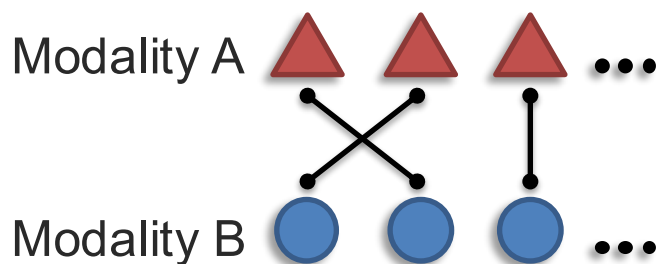
Contextualized representation



Implicit alignment
+ representation

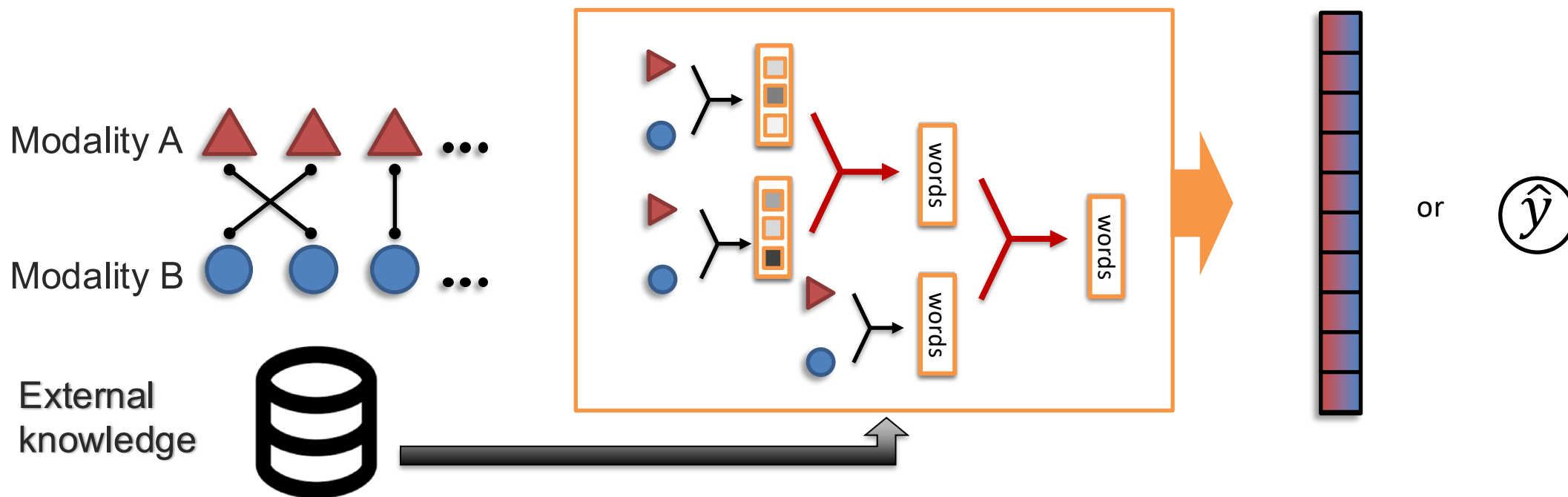
Challenge 3: Reasoning

Definition: Combining knowledge, usually through multiple inferential steps, exploiting multimodal alignment and problem structure.



Challenge 3: Reasoning

Definition: Combining knowledge, usually through multiple inferential steps, exploiting multimodal alignment and problem structure.

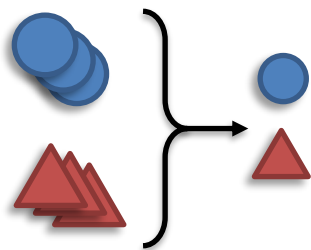


Challenge 4: Generation

Definition: Learning a generative process to produce raw modalities that reflects cross-modal interactions, structure, and coherence.

Sub-challenges:

Summarization



Reduction



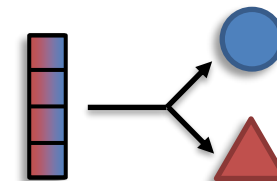
Translation



Maintenance



Creation



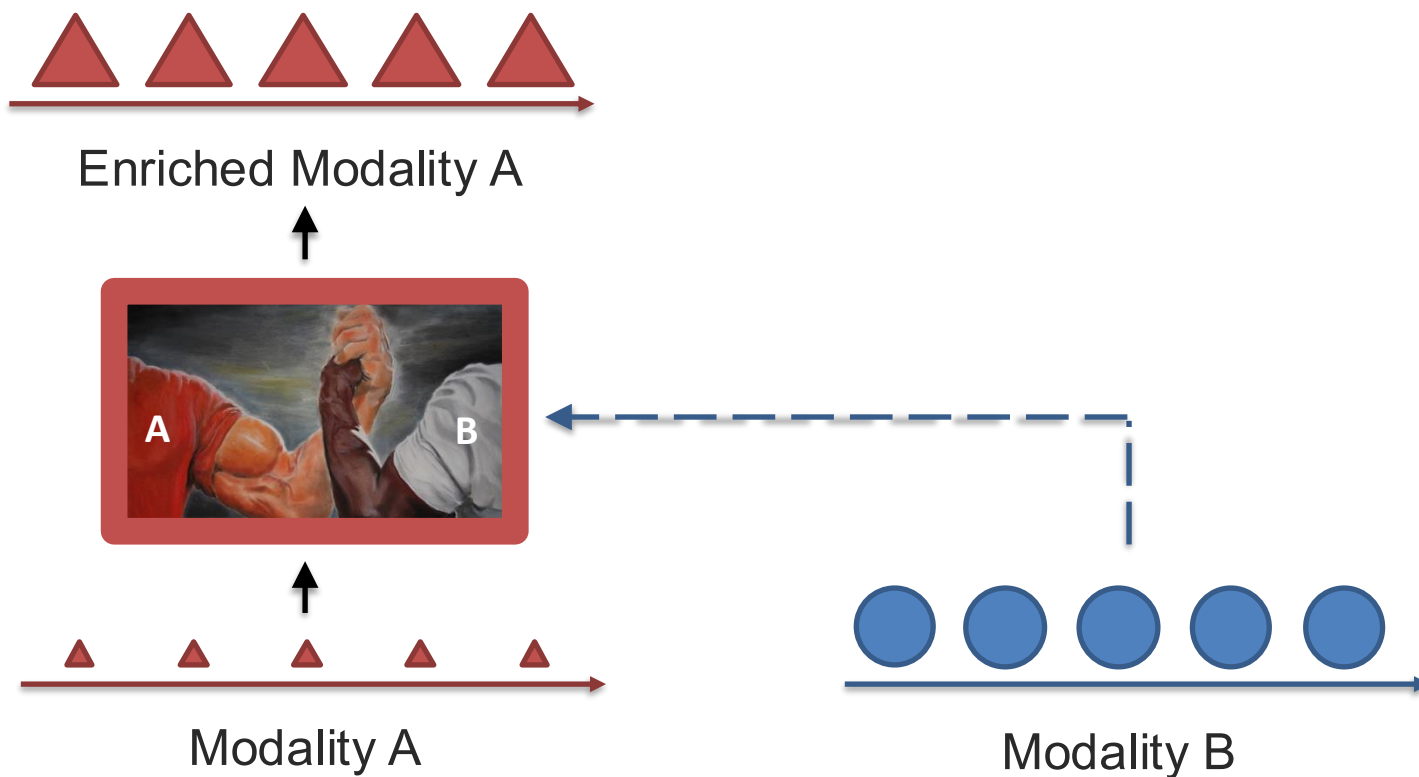
Expansion



Information:
(content)

Challenge 5: Transference

Definition: Transfer knowledge between modalities, usually to help the target modality which may be noisy or with limited resources.

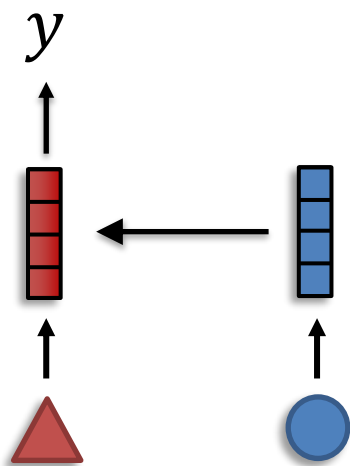


Challenge 5: Transference

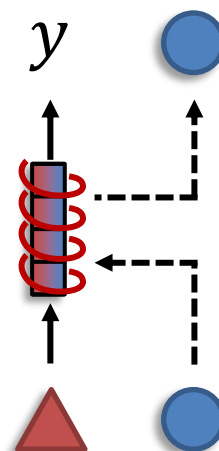
Definition: Transfer knowledge between modalities, usually to help the target modality which may be noisy or with limited resources.

Sub-challenges:

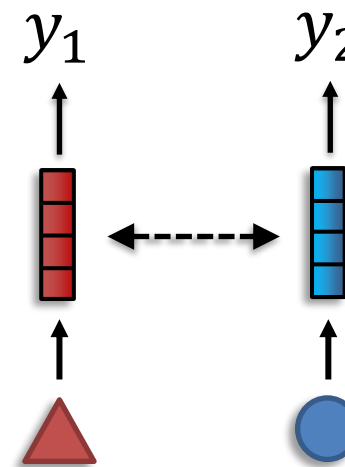
Transfer



Co-learning



Model Induction

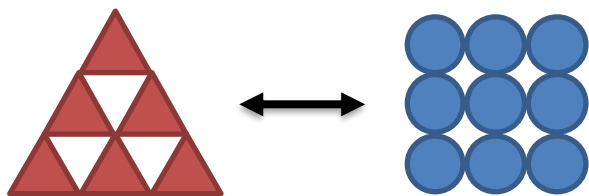


Challenge 6: Quantification

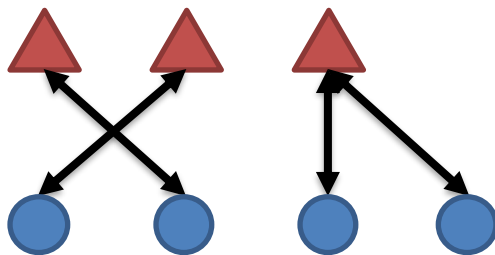
Definition: Empirical and theoretical study to better understand heterogeneity, cross-modal interactions, and the multimodal learning process.

Sub-challenges:

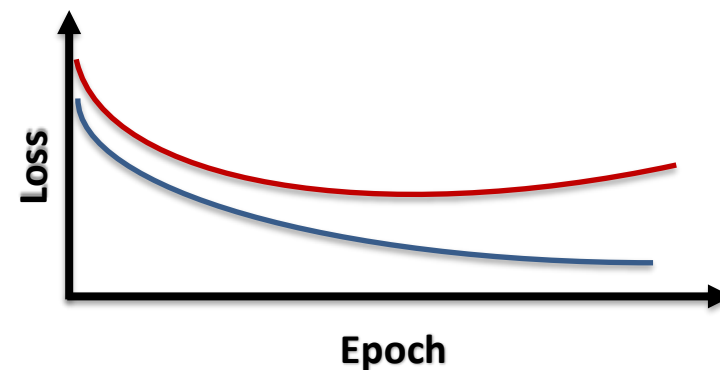
Heterogeneity



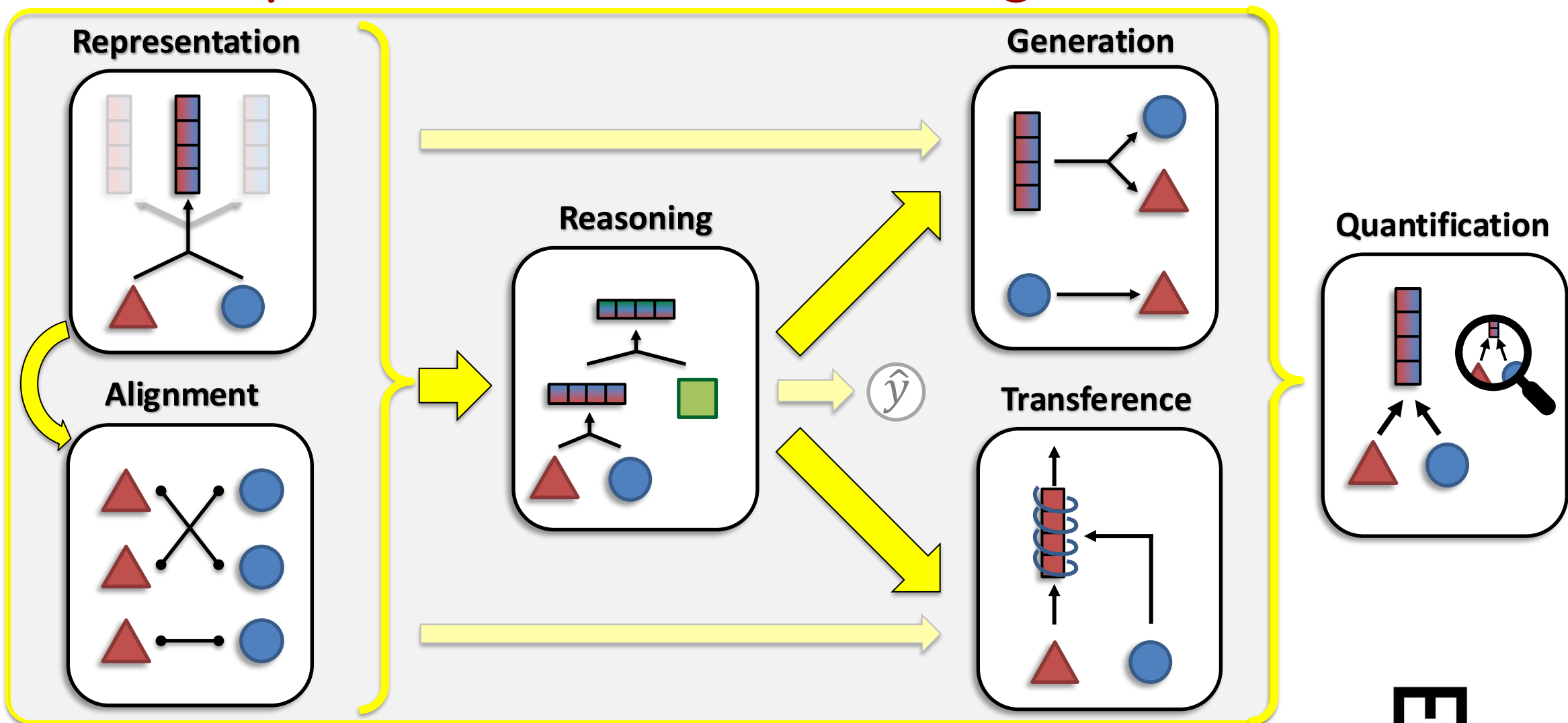
Interactions



Learning



Summary of Core Multimodal Challenges



Native Multimodal Models

- *Native Multimodal Models*: LLMs Trained from scratch with multimodal input (instead of finetuning a trained unimodal LLM)
- Largest public model now: 109B - 2T parameters

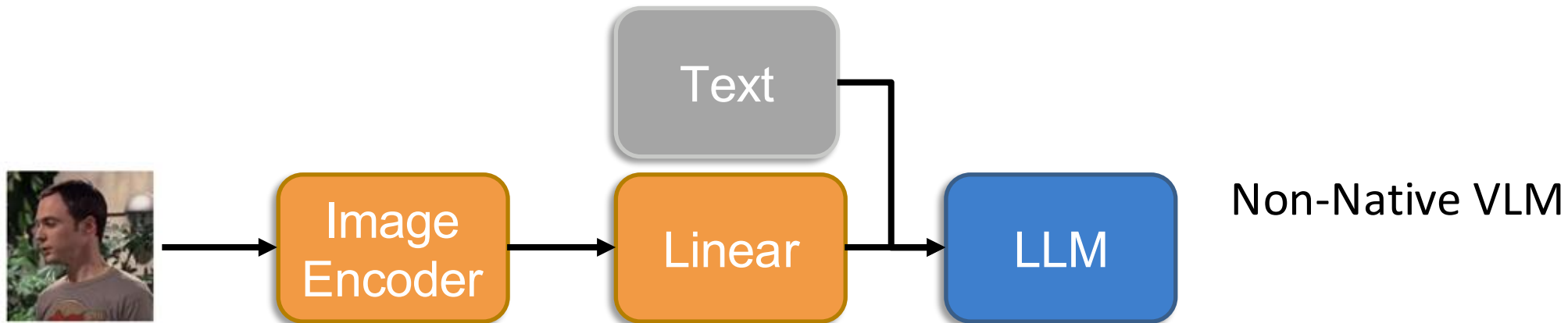
The image shows a promotional graphic for the Llama 4 model family. It features three columns, each representing a different model variant. The background is a light blue gradient. Each model name is in a large, bold, dark font. Below the name, the number of active parameters and total parameters are listed. A short description is provided for each model. At the bottom of each column is a button: 'Preview' for Behemoth, 'Available' for Maverick, and no button for Scout.

Model Name	Active Parameters	Total Parameters	Key Features	Status
Llama 4 Behemoth	288B	2T	The most intelligent teacher model for distillation	Preview
Llama 4 Maverick	17B	400B	Native multimodal with 1M context length	Available
Llama 4 Scout	17B	109B	Industry leading 10M context length Optimized inference	

Native Multimodal Models

- Background

- **Non-native VLMs:** Image encoder paired with frozen trained LLM. The image encoder can either be frozen or trained. Most VLMs now use this structure.

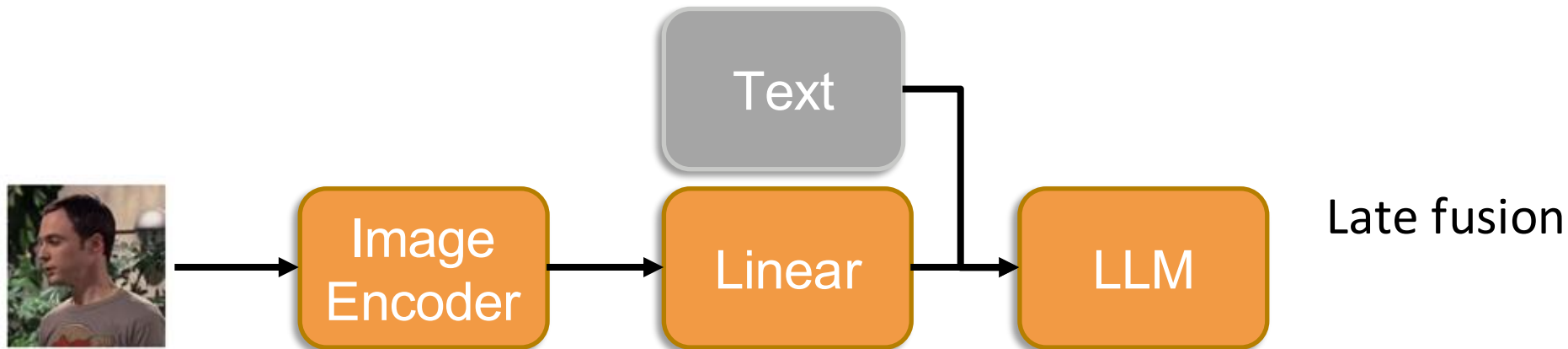


Most current VLMs use this architecture.

Native Multimodal Models

- Background

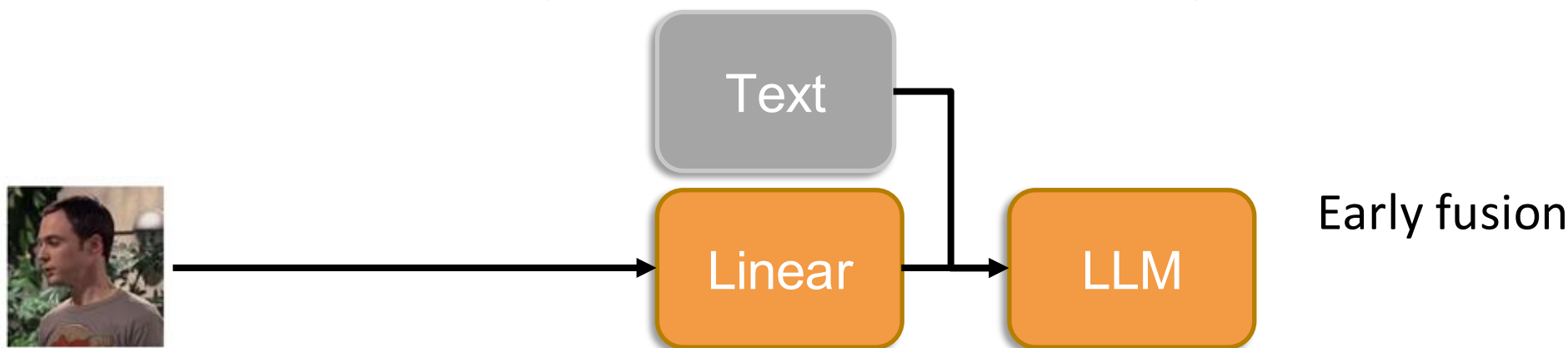
- Non-native VLMs: Image encoder paired with frozen trained LLM. The image encoder can either be frozen or trained. Most VLMs now use this structure.
- Native Multimodal Models: LLMs Trained from scratch with multimodal input
 - **Late fusion:** Image patches -> Image Encoder -> Linear -> LLM.
 - Early fusion: Image patches -> Linear -> LLM (No image encoder!)



Native Multimodal Models

- Background

- Non-native VLMs: Image encoder paired with frozen trained LLM. The image encoder can either be frozen or trained. Most VLMs now use this structure.
- Native Multimodal Models: LLMs Trained from scratch with multimodal input
 - Late fusion: Image patches -> Image Encoder -> Linear -> LLM.
 - **Early fusion**: Image patches -> Linear -> LLM (No image encoder!)

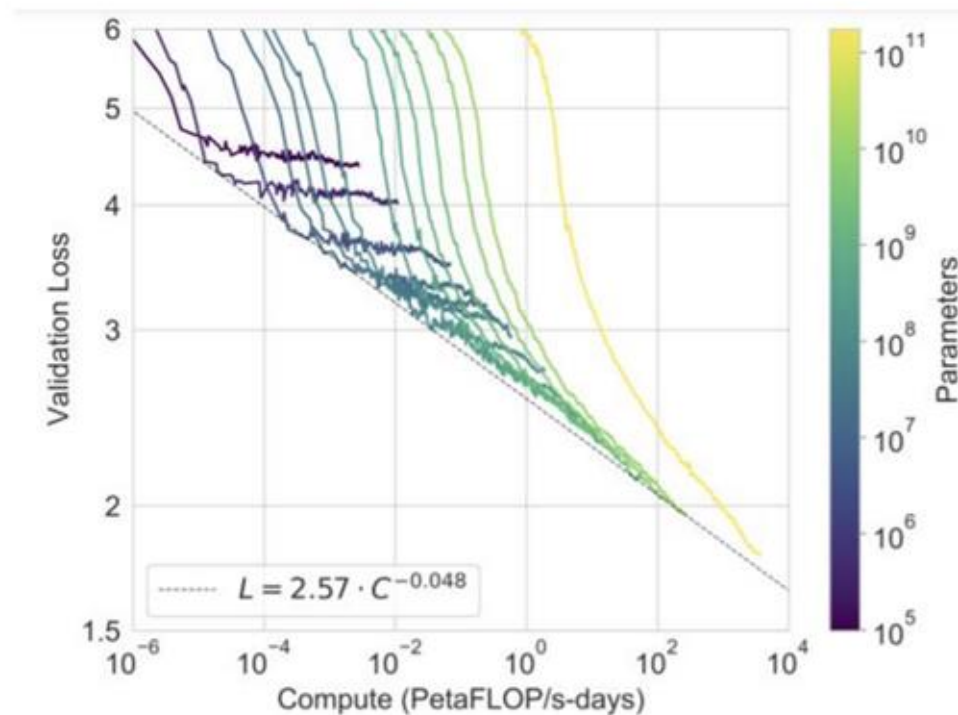


Scaling Laws

- ↯ Bigger model allows models to reach a better performance given sufficient compute
- ↯ Over training models getting popular nowadays

Model	Size (# Parameters)	Training Tokens
LaMDA (Thoppilan et al., 2022)	137 Billion	168 Billion
GPT-3 (Brown et al., 2020)	175 Billion	300 Billion
Jurassic (Lieber et al., 2021)	178 Billion	300 Billion
Gopher (Rae et al., 2021)	280 Billion	300 Billion
MT-NLG 530B (Smith et al., 2022)	530 Billion	270 Billion
<i>Chinchilla</i>	70 Billion	1.4 Trillion

	Training Data	Params	Context length	GQA	Token count	Knowledge cutoff
Llama 3	A new mix of publicly available online data.	8B	8k	Yes	15T+	March, 2023
		70B	8k	Yes		December, 2023



Scaling Laws for Native Multimodal Models

- Studies how multimodal models scale (with model parameters and dataset size) as compared to their unimodal counterparts.

Unimodal scaling:

$$\mathcal{L} \left(\underbrace{N}_{\text{Number of Model Parameters}}, \underbrace{D_j}_{\text{Dataset}} \right) = \underbrace{E_j}_{\text{Minimal Achievable Loss}} + \underbrace{\frac{A_j}{N^{\alpha_j}}}_{\text{Functional Approximation Error}} + \underbrace{\frac{B_j}{|D_j|^{\beta_j}}}_{\text{Convergence Error}} \quad (2)$$

For Modality j

Multimodal scaling:

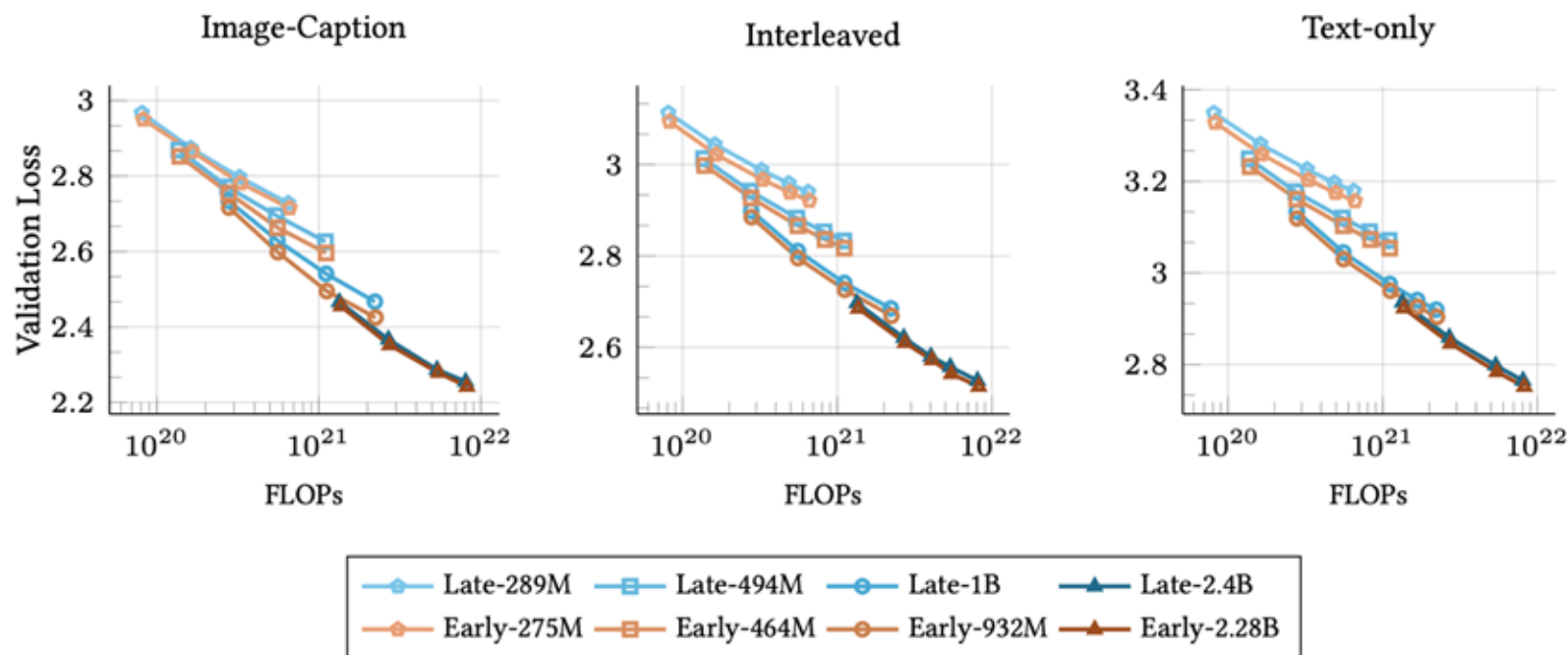
$$\mathcal{L} (N, D_i, D_j) = \left[\frac{\mathcal{L} (N, D_i) + \mathcal{L} (N, D_j)}{2} \right] - \underbrace{C_{i,j}}_{\text{Competition in Functional Approximation}} + \underbrace{\frac{A_{i,j}}{N^{\alpha_{i,j}}}}_{\text{Competition in Optimization Process}} + \frac{B_{i,j}}{|D_i| + |D_j|^{\beta_{i,j}}} \quad (4)$$

Loss if Datasets Were Modeled Independently

Maximum Level of Synergy

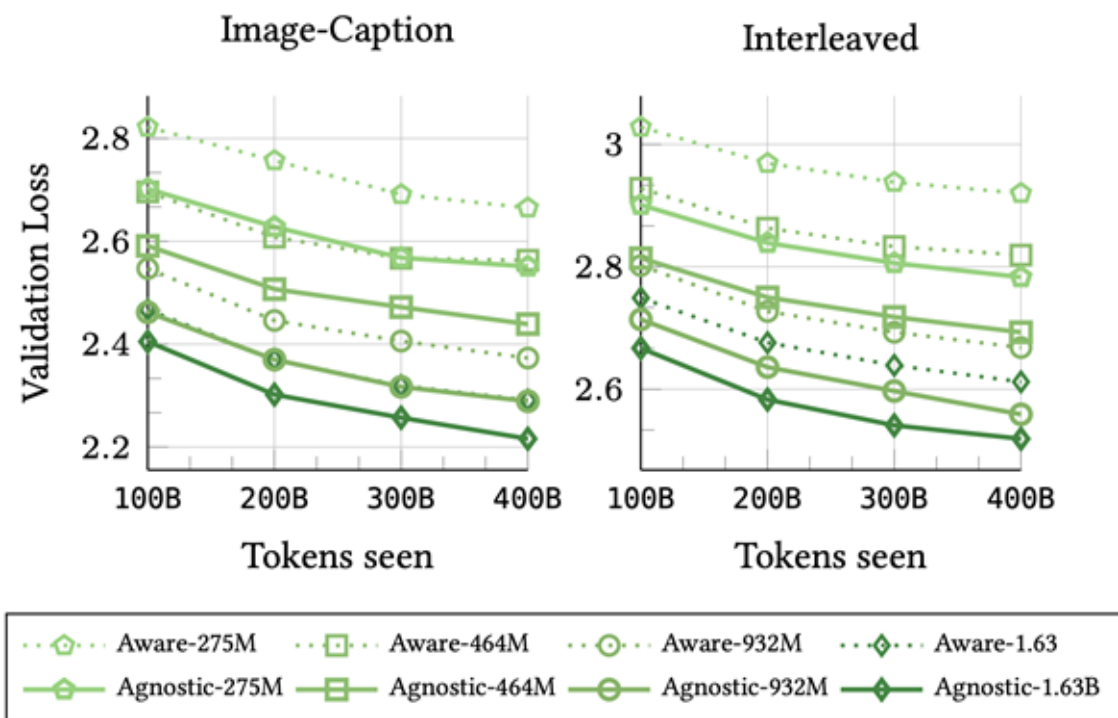
Scaling Laws for Native Multimodal Models

- Early fusion models hold small advantage on small scales.
- On larger scales, **both architectures perform similarly**. (We don't actually need image encoders!)
- **NMMs scale similarly to unimodal LLMs**, with slightly varying scaling exponents depending on the target data type and training mixture



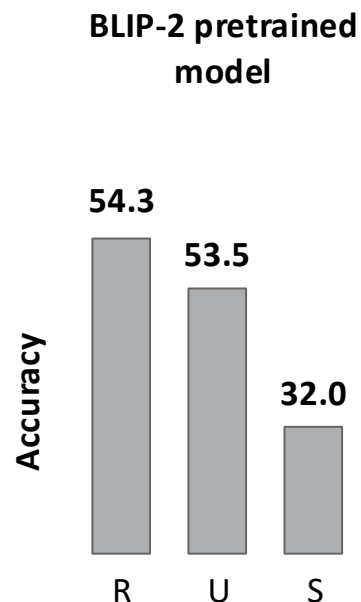
Scaling Laws for Native Multimodal Models

- Sparse structure like **MOE significantly benefits NMMs** at the same inference cost
- In an MOE structure, Modality-aware design (having separate image/text experts) performs **worse** than modality-agnostic design (unified experts for both image/text tokens)



One model for everything?

Video sarcasm detection

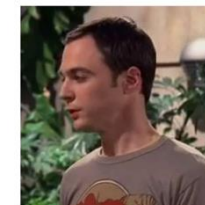


Y: Sarcasm

X_ℓ : Spoken language

It's just a privilege to watch your mind at work.

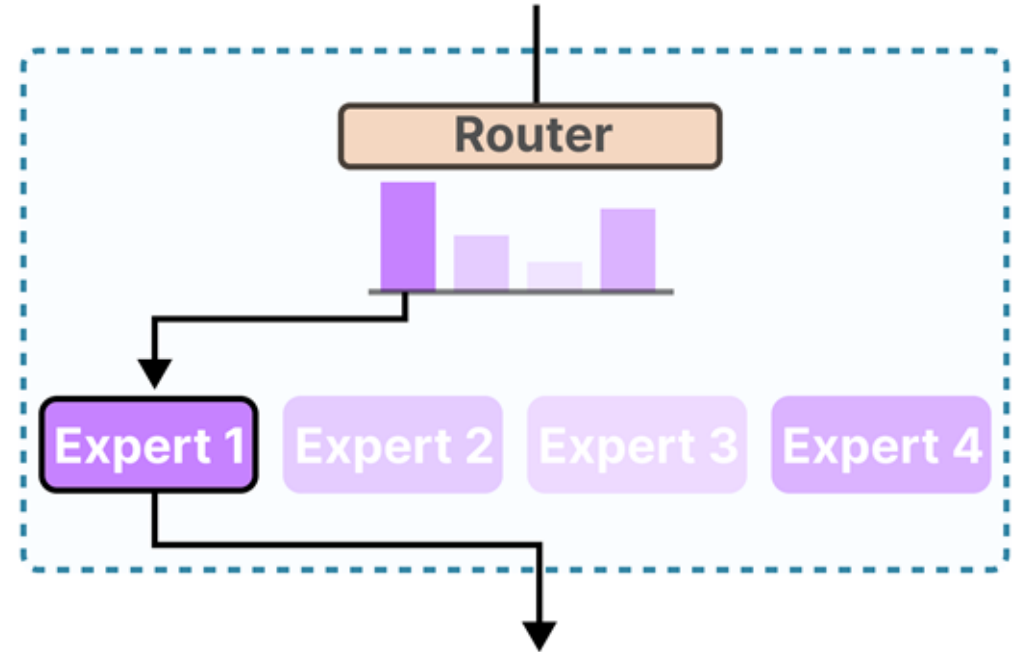
X_{av} : Audio + visual



Neutral tone + straight face

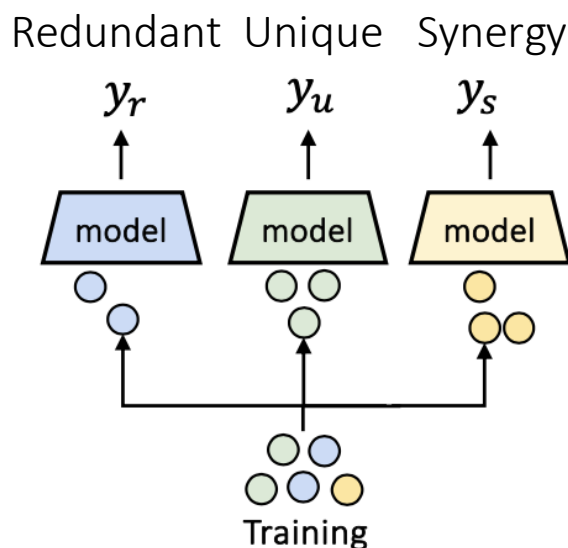
Efficient Training: Mixture of Experts

- ↵ Train multiple parallel networks (experts) simultaneously
- ↵ During each forward pass, only activate k experts
- ↵ Saves compute & GPU memory
- ↵ Deepseek R1: 671B, only 37B activated, performance on par with OpenAI o1-mini

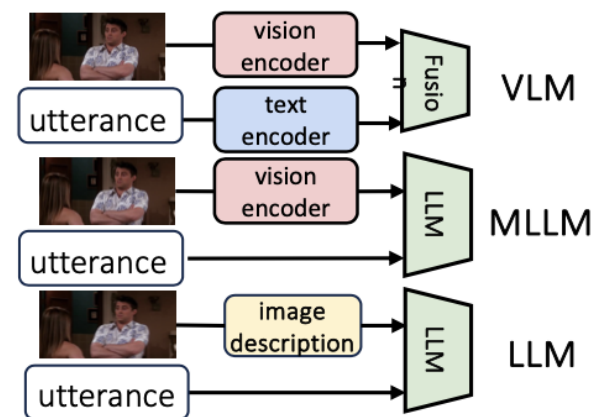
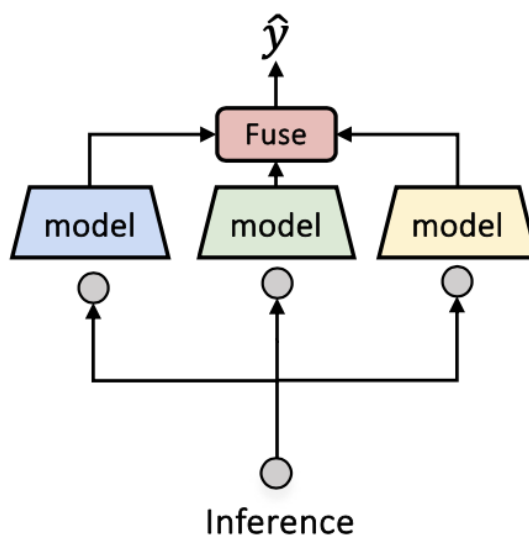


Mixture of Multimodal Interaction Experts

One model for everything -> specialized models for each interaction



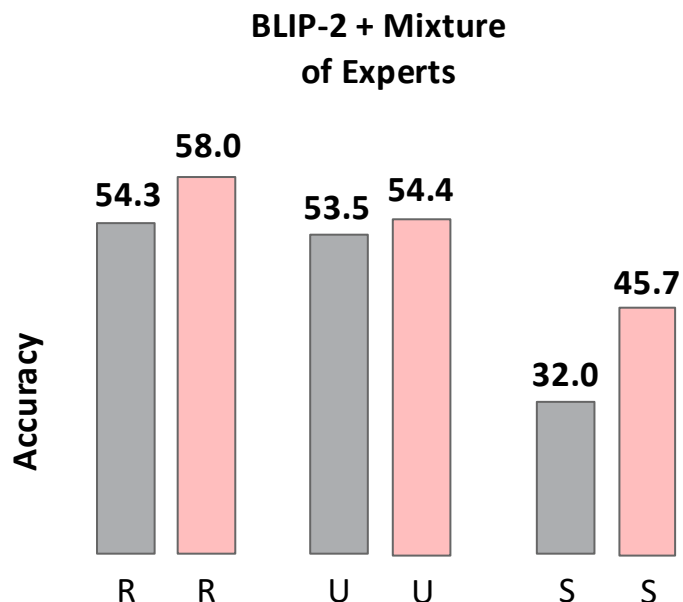
Datapoint-level
quantification
of interactions



Mixture of Multimodal Interaction Experts

One model for everything -> specialized models for each interaction

Video sarcasm detection



The car is as fast as a cheetah.



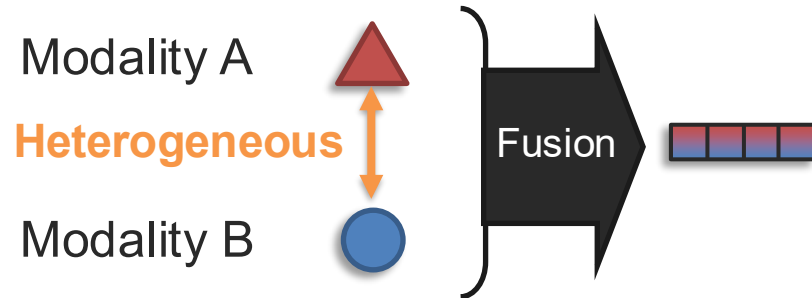
[Yosef et al., EMNLP 23]



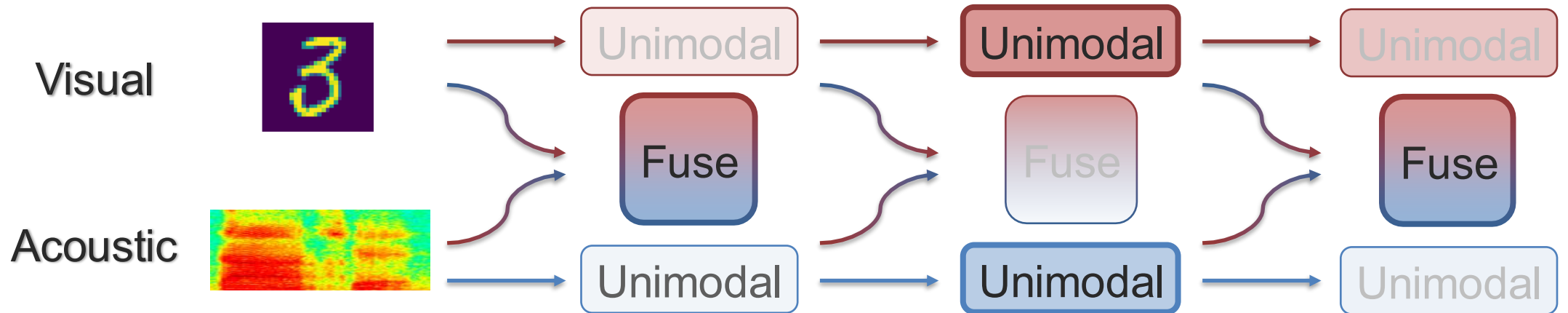
Can you please pass the cow?

[Hessel et al., ACL 23]

Dynamic Early Fusion



Idea: Deciding when to fuse in early fusion



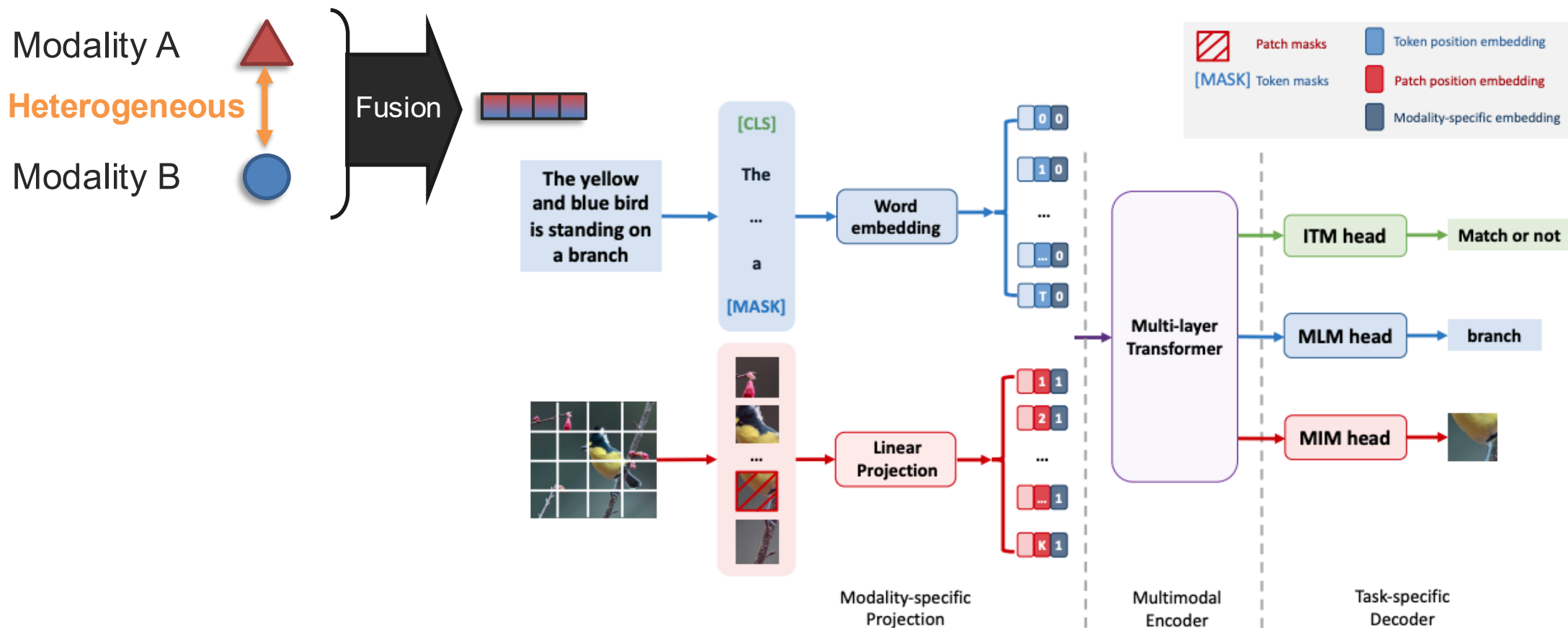
[Xue and Marculescu, Dynamic Multimodal Fusion, arxiv 2022]

[Xu et al., MUFASA: Multimodal Fusion Architecture Search for Electronic Health Records. AAAI 2021]

[Liu et al., DARTS: Differentiable Architecture Search. ICLR 2019]

Fusion with Heterogeneous Modalities

Example: From feature fusion to early fusion



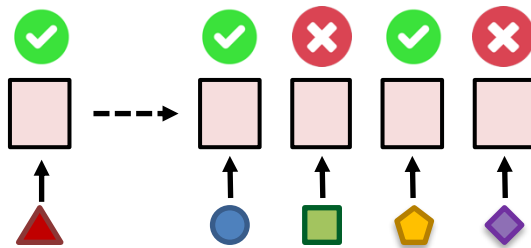
[Liang et al., High-modality Multimodal Transformer. TMLR 2022]

[Gui et al., Training Vision-Language Transformers from Captions. arxiv 2022]

Fusion with Heterogeneous Modalities

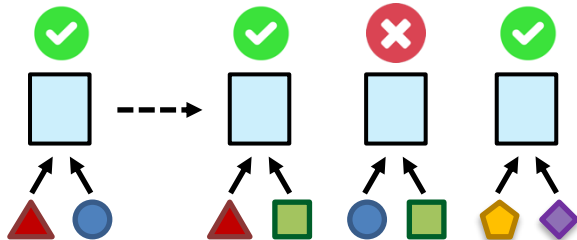
Information transfer, transfer learning perspective

1a. Estimate modality heterogeneity via transfer



(Implicitly captures heterogeneity)

1b. Estimate interaction heterogeneity via transfer



2a. Compute modality heterogeneity matrix

	0				
	1	0			
	3	2	0		
	1	2	3	0	
	5	4	6	3	0

3. Determine parameter clustering

$$U_1 = \{U_1, U_2, U_4\}$$

$$U_2 = \{U_3\}$$

$$U_3 = \{U_5\}$$

$$C_1 = \{C_{12}, C_{13}, C_{45}\}$$

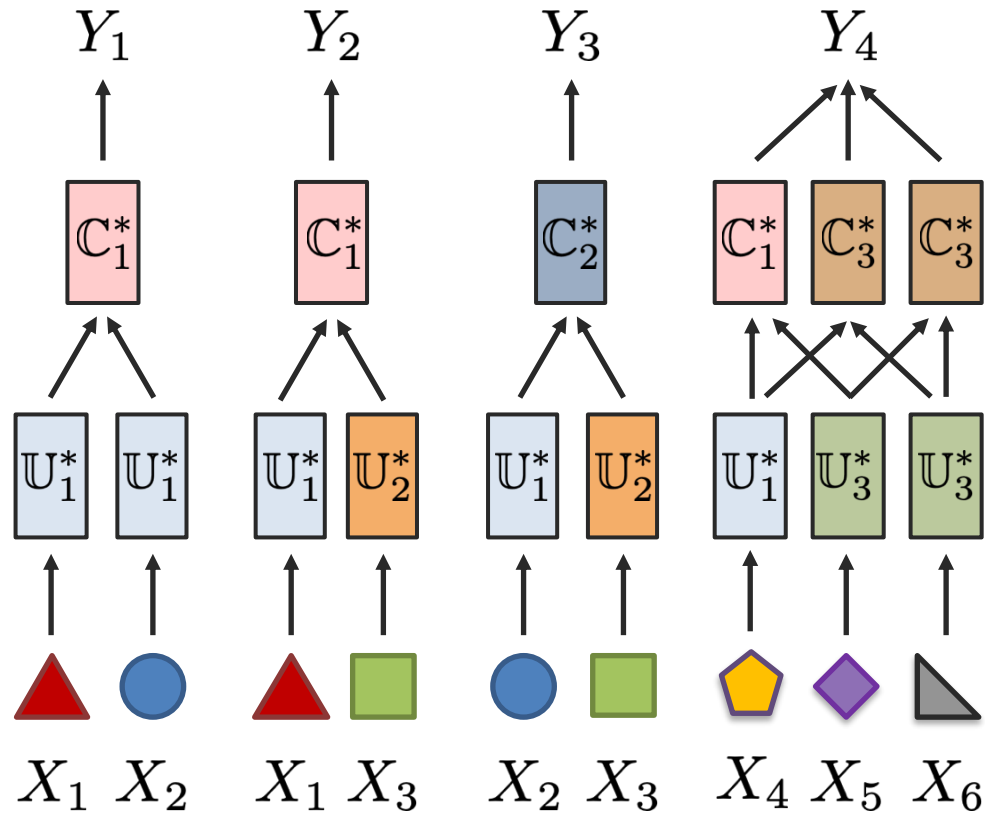
$$C_2 = \{C_{23}\}$$

2b. Compute interaction heterogeneity matrix

	0			
	1	0		
	3	2	0	
	1	2	4	0

Fusion with Heterogeneous Modalities

Information transfer, transfer learning perspective



Nonlinear Fusion

Kinetics dataset



(a) headbanging



(c) shaking hands



(e) robot dancing



(g) riding a bike



Adding more modalities should always help?

Modalities: **RGB** (video clips)

A (Audio features)

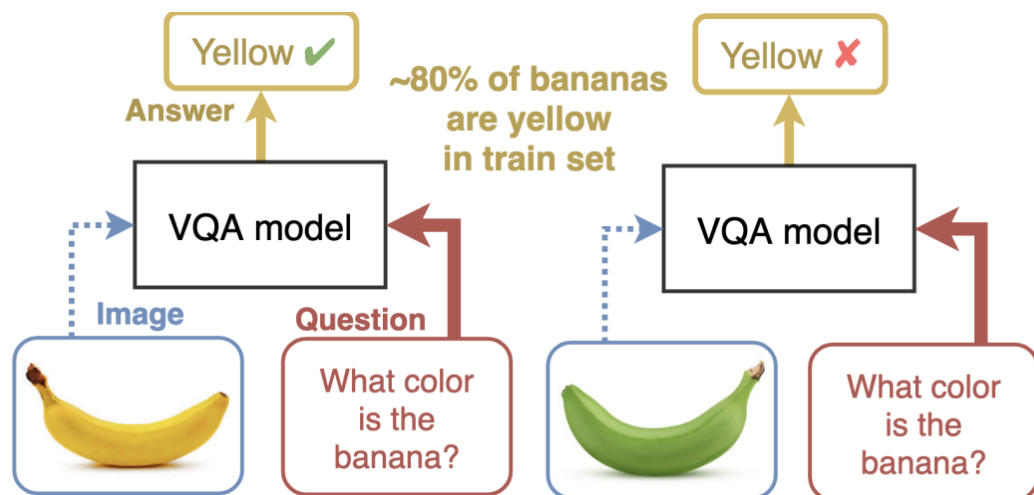
OF (optical flow - motion)

Dataset	Multi-modal	V@1	Best Uni	V@1	Drop
Kinetics	A + RGB	71.4	RGB	72.6	-1.2
	RGB + OF	71.3	RGB	72.6	-1.3
	A + OF	58.3	OF	62.1	-3.8
	A + RGB + OF	70.0	RGB	72.6	-2.6

But sometimes multimodal doesn't help! **Why?**

Unimodal Biases

Finding: VQA models answer the question without looking at the image

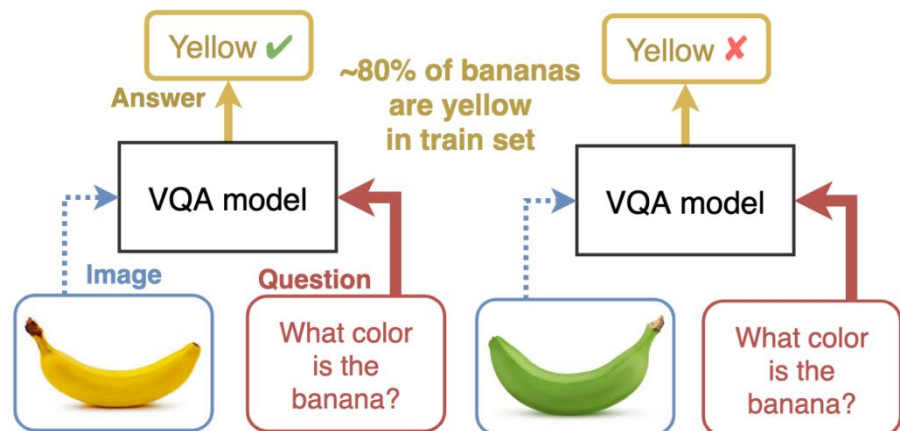


Finding: Image captioning models capture spurious correlations between gender and generated actions.



Unimodal Biases

VQA models answer the question without looking at the image

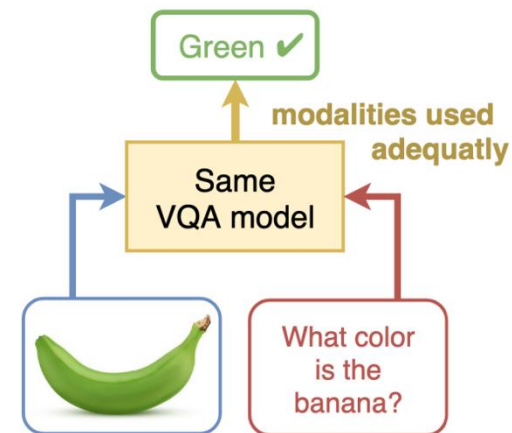


Balancing modalities

Balancing training



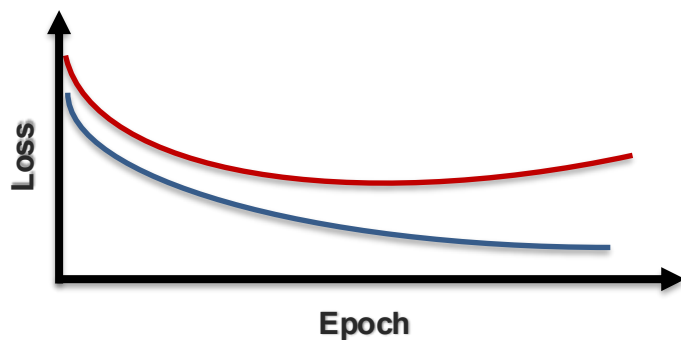
Not the case when trained with RUBi



Optimization Challenges

2 explanations for drop in performance:

1. Multimodal networks are more prone to overfitting due to **increased complexity**
2. Different modalities overfit and generalize at **different rates**



Key idea 1: compute overfitting-to-generalization ratio (OGR)

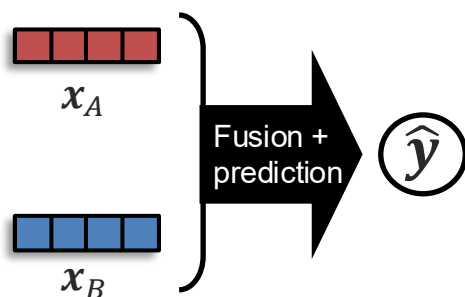


Gap between training and valid loss

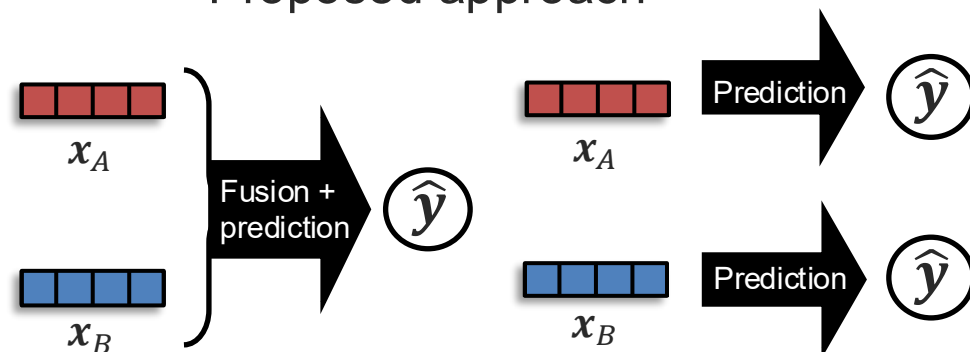
OGR wrt each modality tells us
how much to train that modality

Optimization Challenges

Conventional approach



Proposed approach



Key idea 2: Simultaneously train unimodal networks to estimate OGR wrt each modality



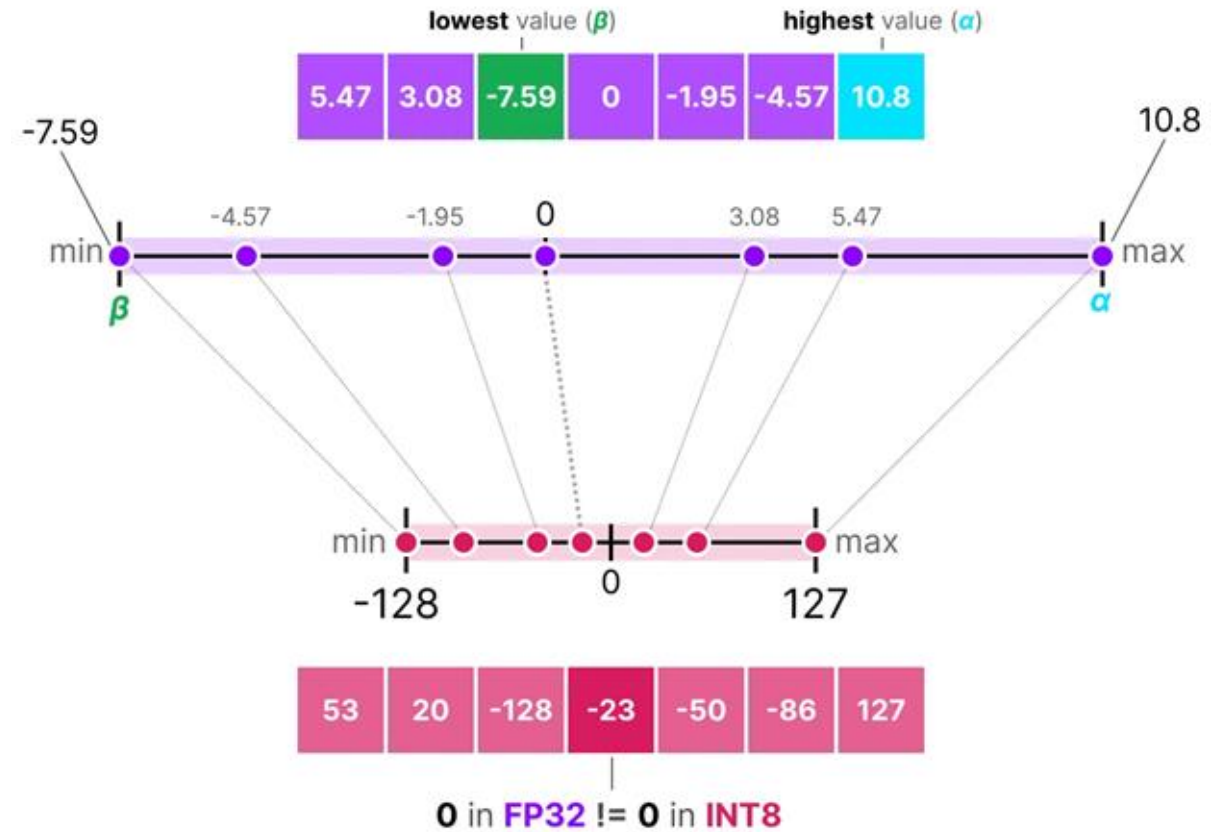
Reweight multimodal loss using unimodal OGR values



Allows to better balance generalization & overfitting rate of different modalities

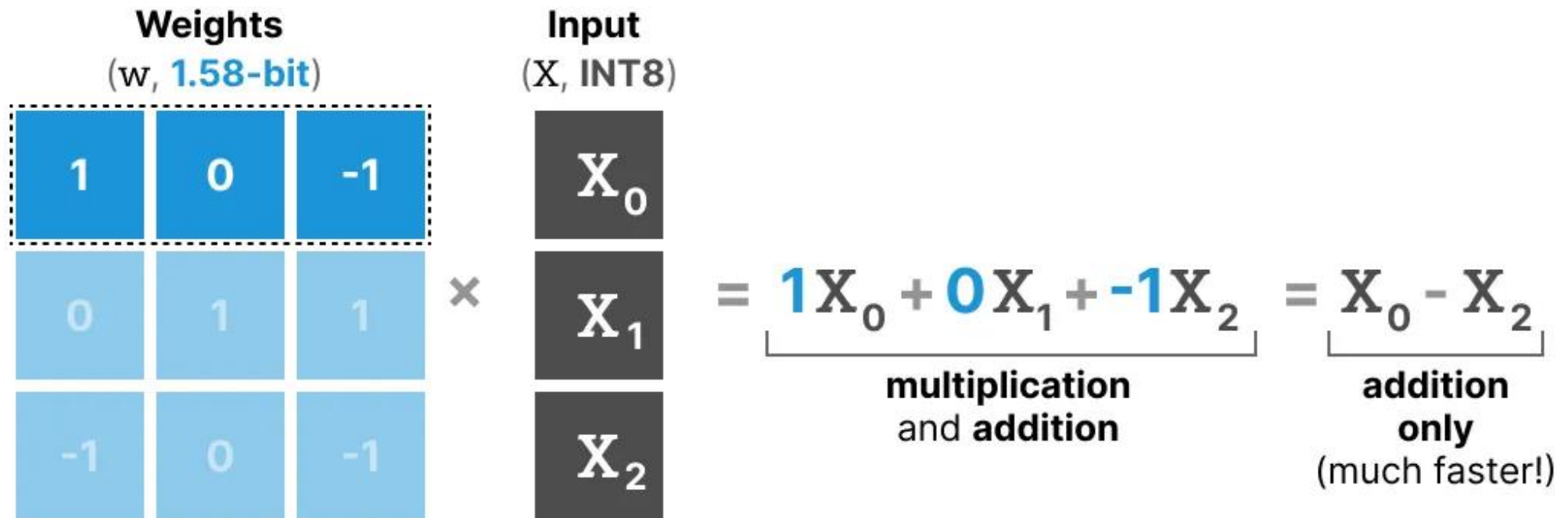
Efficient Inference: Quantization

- ↯ Range Clipping
- ↯ Scale & Shift
- ↯ Convert to lower bits
- ↯ Calibration



Efficient Inference: Quantization

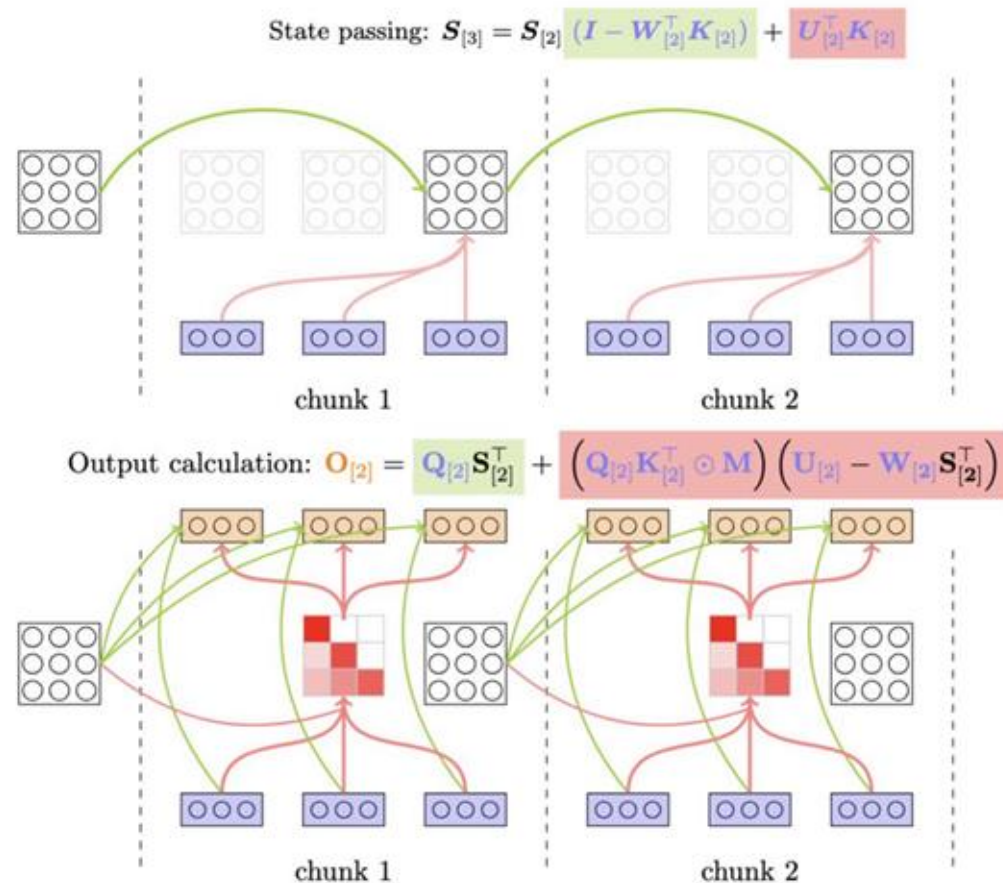
- Quantize weights to ternary $\{-1, 0, 1\}$
- New hardware needed to be truly 1.58 bits / $\log_2(3)$



Linear Attention

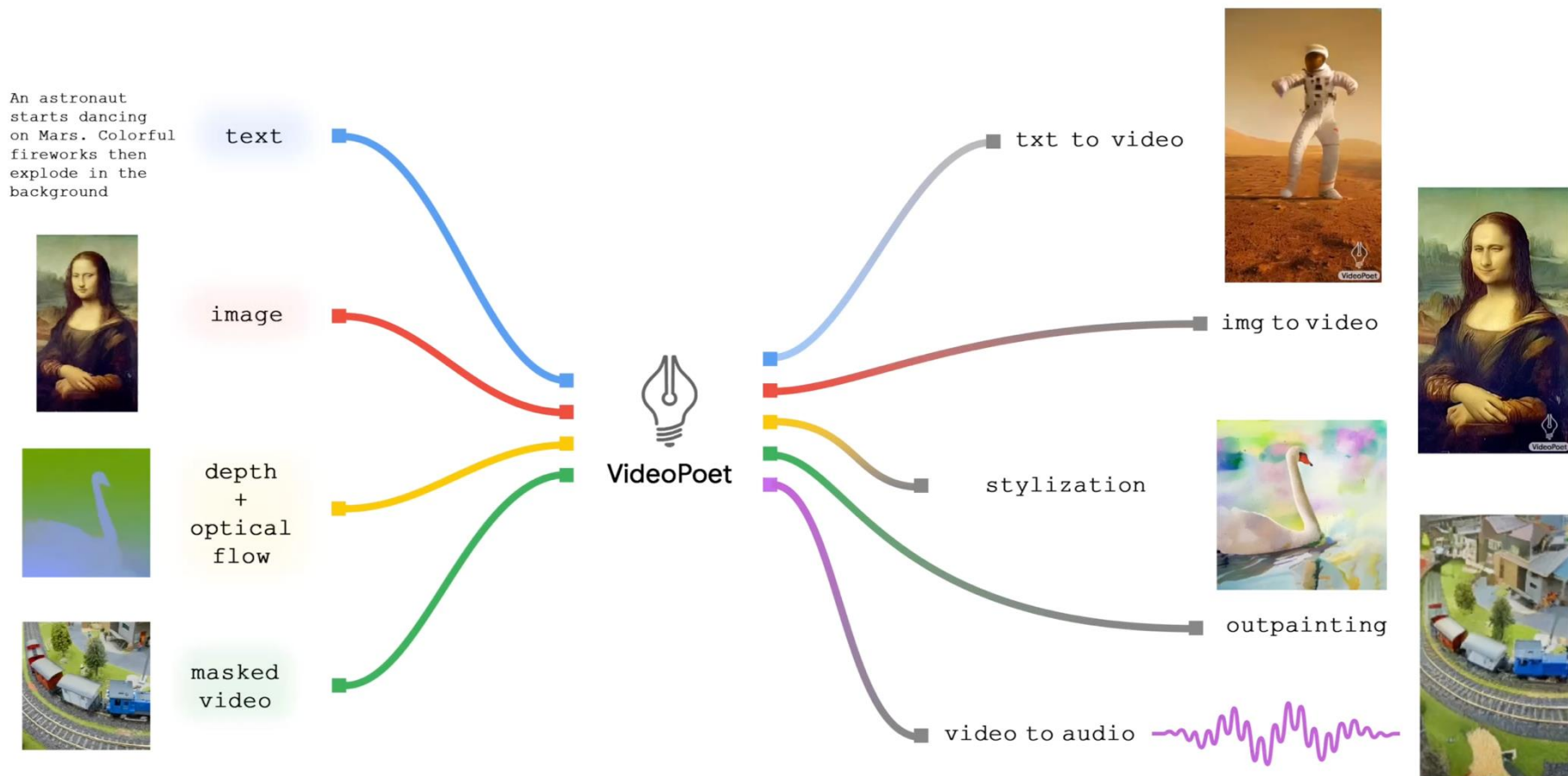
- Time Complexity of attention scales at $O(n^2)$
- Solution: Linear Attention
- Key idea: Combine Transformer Style Full Attention with RNN
- Split sentence into chunks, run attention within chunk, then combine them in RNN style aggregation
- Promising performance on par with full attention

$$h = \text{softmax} \left(\frac{XW_q W_k^T X^T}{\sqrt{d}} \right) XW_v$$



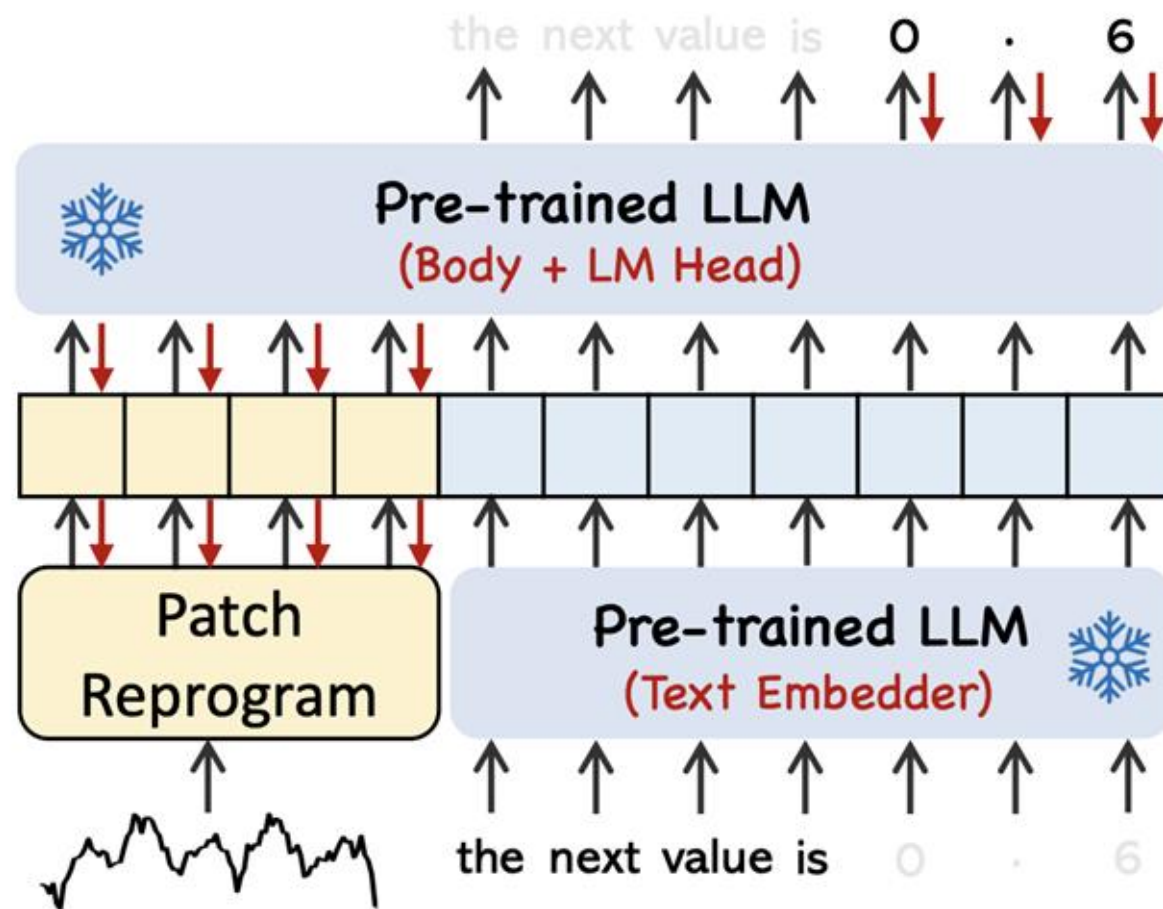
Any-to-any Models

Multimedia, content creation, creativity and the arts



Time-series Multimodal Models

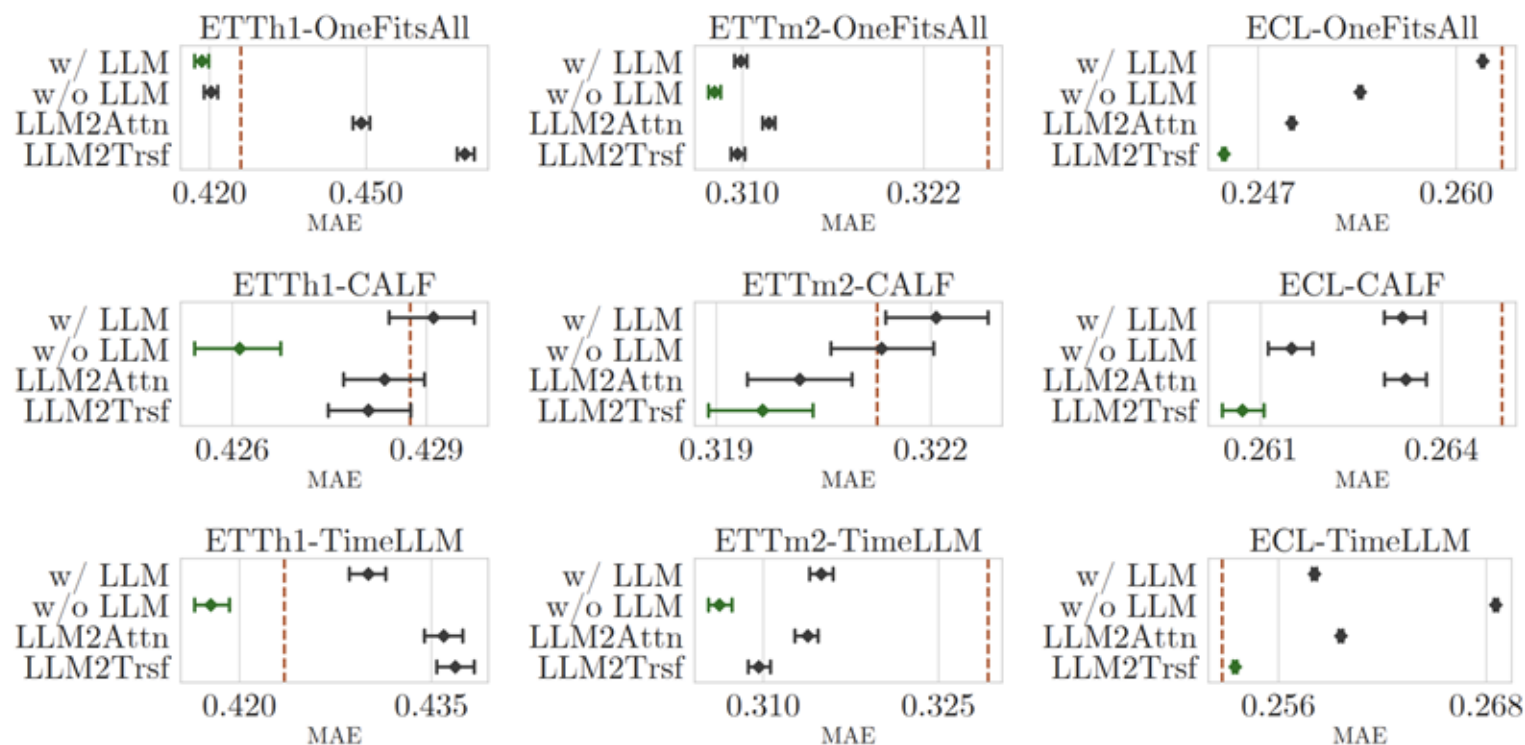
- Typically trained & aligned the same way as vision language models (alignment + instruction tuning)
- Works for both analysis and prediction
- Example: Time-LLM, OneFitsAll



Time-series Multimodal Models

- But some current time series LLMs have questionable performance. Replacing LLM with a simple attention layer doesn't significantly degrade performance (sometimes even better).

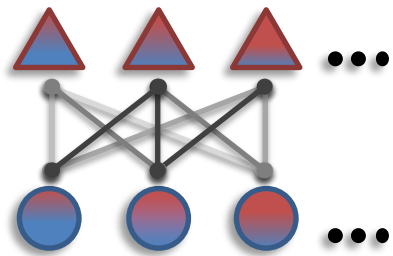
* Lower is better



Multimodal Reasoning Agents

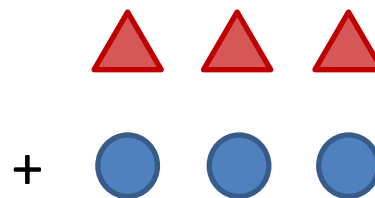
Solving hard problems using step-by-step reasoning and action taking in multiple modalities

It's just a privilege to watch your mind at work.



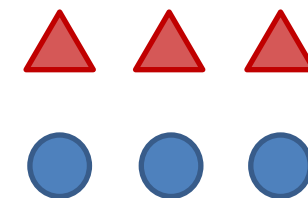
Multimodal representation

*This person is being sarcastic.
They seem to be close friends.*



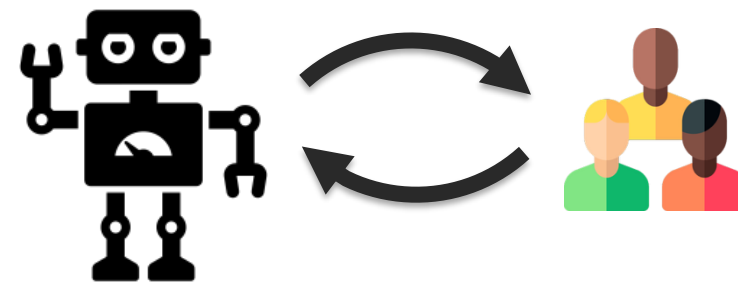
*(quote previous episodes)
(highlight multimodal information)*

*Here's a story of them in
a different culture...*



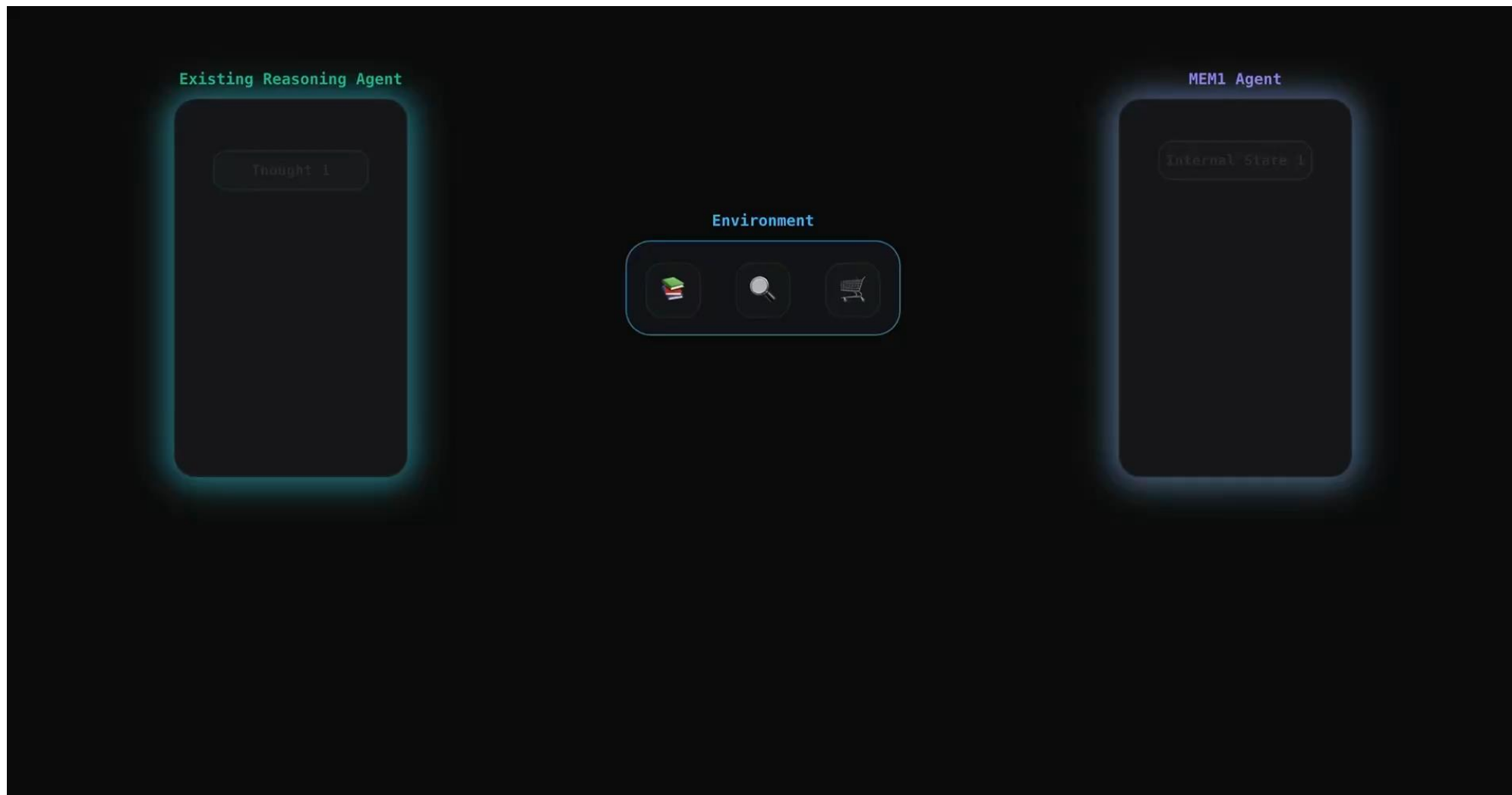
(generate future episodes)

Reasoning context is getting longer and longer...



MEM1: Memory Efficient Reasoning

<https://github.com/MIT-MI/MEM1>



Self-Evolving Agents

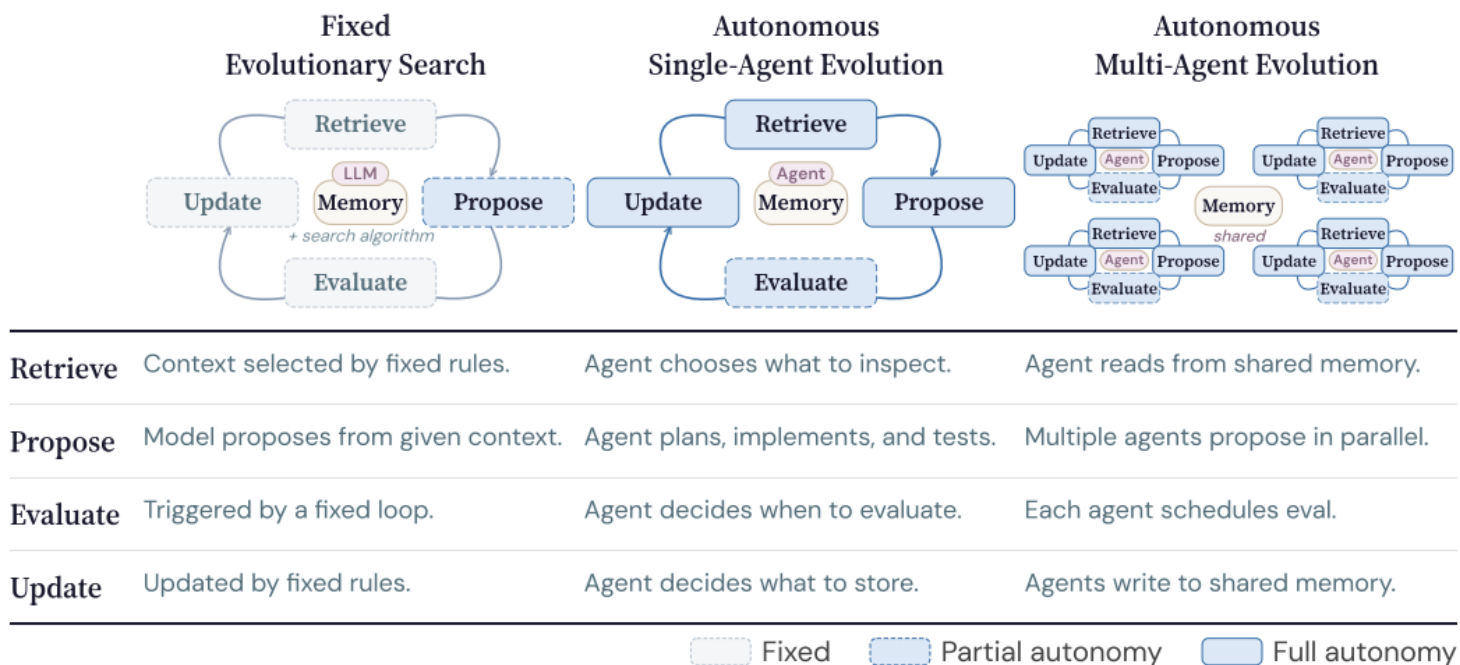
MEM1 is an agent that learns to manage its memory. Other agents manage tools, APIs, search, coding etc.

Question 1: How can agents automatically decide what functions to use for a given task?

Answer: Self-evolving agents

Question 2: How can agents automatically delegate tasks to different agents and organize themselves?

Answer: Self-evolving multi-agent systems



Self-Evolving Agents

Agents are now helping with scientific discovery

nature

Explore content ▾ About the journal ▾ Publish with us ▾

[nature](#) > [articles](#) > [article](#)

Article | [Open access](#) | Published: 14 December 2023

Mathematical discoveries from program search with large language models

[Bernardino Romera-Paredes](#) , [Mohammadamin Barekatin](#), [Alexander Novikov](#), [Matej Balog](#), [M. Pawan Kumar](#), [Emilien Dupont](#), [Francisco J. R. Ruiz](#), [Jordan S. Ellenberg](#), [Pengming Wang](#), [Omar Fawzi](#), [Pushmeet Kohli](#)  & [Alhussein Fawzi](#) 

[Nature](#) 625, 468–475 (2024) | [Cite this article](#)

321k Accesses | 423 Citations | 1054 Altmetric | [Metrics](#)

May 14, 2025 Science

AlphaEvolve: A Gemini-powered coding agent for designing advanced algorithms

AlphaEvolve team

Share 

AlphaEvolve and OpenEvolve

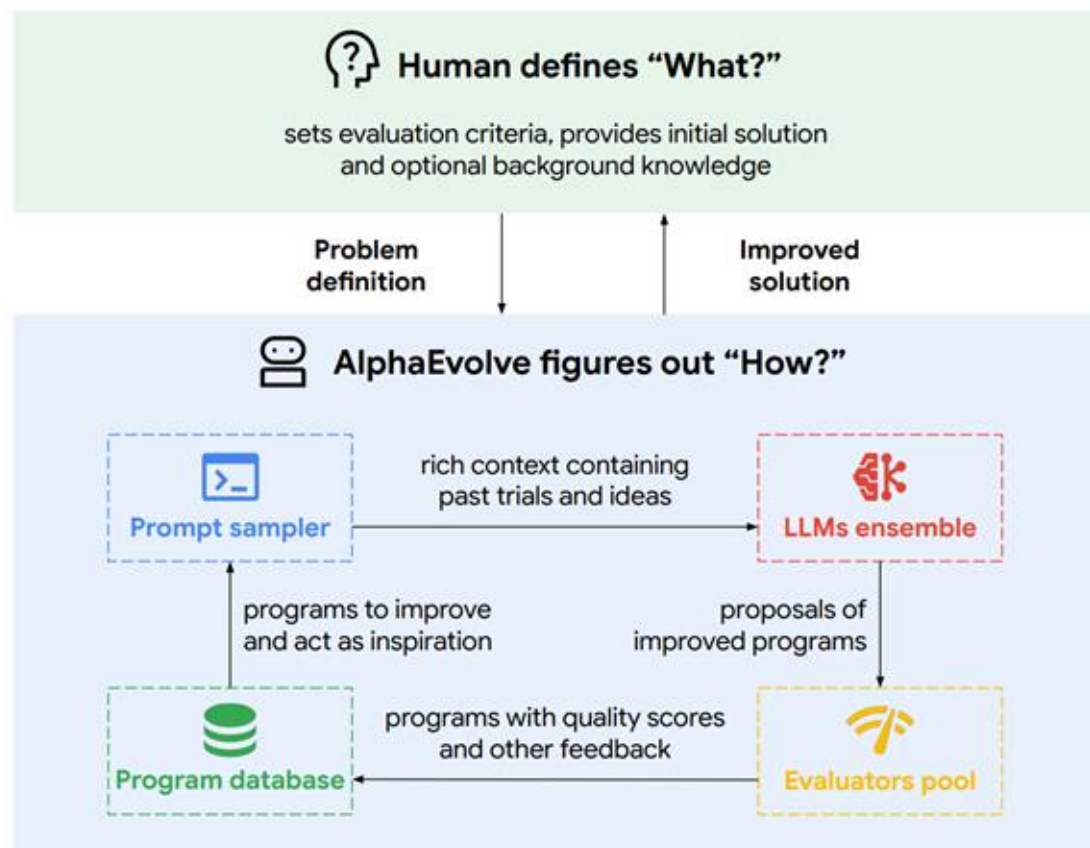


Figure 1 | AlphaEvolve high-level overview.

```
# EVOLVE-BLOCK START
"""Image classification experiment in jaxline."""

import jax
...
# EVOLVE-BLOCK-END

...

# EVOLVE-BLOCK-START
class ConvNet(hk.Module):
    def __init__(self, num_classes): ...
    def __call__(self, inputs, is_training): ...

def sweep():
    return hyper.zipit(...)
# EVOLVE-BLOCK-END

...

def evaluate(eval_inputs) -> dict(str, float):
    ...
    return metrics
```

(a) ?

The current model uses a simple ResNet architecture with only three ResNet blocks. We can improve its performance by increasing the model capacity and adding regularization. This will allow the model to learn more complex features and generalize better to unseen data. We also add weight decay to the optimizer to further regularize the model and prevent overfitting. AdamW is generally a better choice than Adam, especially with weight decay.

```
<<<<<< SEARCH
self_block1 = ResNetBlock(num_channels)
self_block2 = ResNetBlock(num_channels = 2, stride=2)
self_block3 = ResNetBlock(num_channels = 4, stride=2)
=====
self_block1 = ResNetBlock(num_channels)
self_block2 = ResNetBlock(num_channels, stride=1)
self_block3 = ResNetBlock(num_channels = 2, stride=2)
self_block4 = ResNetBlock(num_channels = 2, stride=1)
self_block5 = ResNetBlock(num_channels = 4, stride=2)
self_block6 = ResNetBlock(num_channels = 4, stride=1)
>>>>>> REPLACE
<<<<<< SEARCH
def optimizer(self, learning_rate):
    return optax.adam(learning_rate)
=====
def optimizer(self, learning_rate):
    return optax.adam(learning_rate, weight_decay=1e-4)
>>>>>> REPLACE
```

(c) ?

Act as an expert software developer. Your task is to iteratively improve the provided codebase. [...]

- Prior programs

Previously we found that the following programs performed well on the task at hand:

top_1_acc: 0.796; neg_eval_log_loss: 0.230; average_score: 0.513

"""Image classification experiment in jaxline."""

[...]

class ConvNet(hk.Module):

"""Network"""

```
def __init__(self, num_channels=32, num_output_classes=10):
    super().__init__()
    self_conv1 = hk.Conv2D(num_channels, kernel_shape=3)
    self_conv2 = hk.Conv2D(num_channels + 2, kernel_shape=3)
    self_conv3 = hk.Conv2D(num_channels + 4, kernel_shape=3)
    self_logits_module = hk.Linear(num_output_classes)
[...]
```

- Current program

Here is the current program we are trying to improve (you will need to propose a modification to it below):

top_1_acc: 0.862; neg_eval_log_loss: 0.387; average_score: 0.624

"""Image classification experiment in jaxline."""

[...]

class ConvNet(hk.Module):

"""Network"""

```
def __init__(self, num_channels=32, num_output_classes=10):
    super().__init__()
    self_conv1 = hk.Conv2D(num_channels, kernel_shape=3)
    self_block1 = ResNetBlock(num_channels)
    self_block2 = ResNetBlock(num_channels + 2, stride=2)
    self_block3 = ResNetBlock(num_channels + 4, stride=2)
    self_logits_module = hk.Linear(num_output_classes)
[...]
```

SEARCH/REPLACE block rules:

[...]

Make sure that the changes you propose are consistent with each other. For example, if you refer to a new config variable somewhere, you should also propose a change to add that variable.

Example:

[...]

Task

Suggest a new idea to improve the code that is inspired by your expert knowledge of optimization and machine learning.

Describe each change with a SEARCH/REPLACE block.

(b) >

Problems Solved by AlphaEvolve

$$\max_{-1/2 \leq s \leq 1/2} \int_{\mathbb{R}} f(t-x)f(x) dx \geq \mathbb{C} \left(\int_{-1/4}^{1/4} f(x) dx \right)^2$$

1.5098 → 1.5053

$$\|f * f\|_2^2 \leq \mathbb{C} \|f * f\|_1 \|f * f\|_\infty$$

0.8892 → 0.8962

$$\max_{-1/2 \leq s \leq 1/2} \left| \int_{\mathbb{R}} f(t-x)f(x) dx \right| \geq \mathbb{C} \left(\int_{-1/4}^{1/4} f(x) dx \right)^2$$

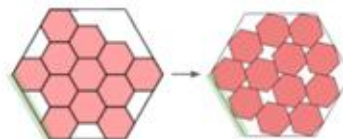
1.4581 → 1.4557

$$A(f)A(\hat{f}) \geq \mathbb{C}'''$$

0.3523 → 0.3521

Analysis

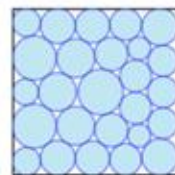
Hexagon outer edge
4.000 → 3.942



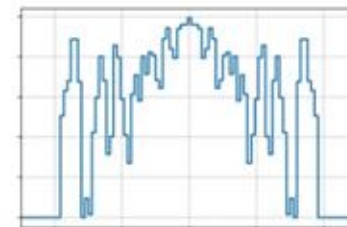
Max distance/min distance
12.890 → 12.889



Sum of radii
2.6340 → 2.6358



Geometry



$$\sup_{x \in [-2,2]} \int_{-1}^1 f(t)g(x+t) dt \geq \mathbb{C}$$

0.380926 → 0.380924

$$|A+B| \ll |A|$$

$$|A-B| \gg |A|^{\mathbb{C}}$$

1.1446 → 1.1584

Combinatorics

Evolve Solutions via Tree Search

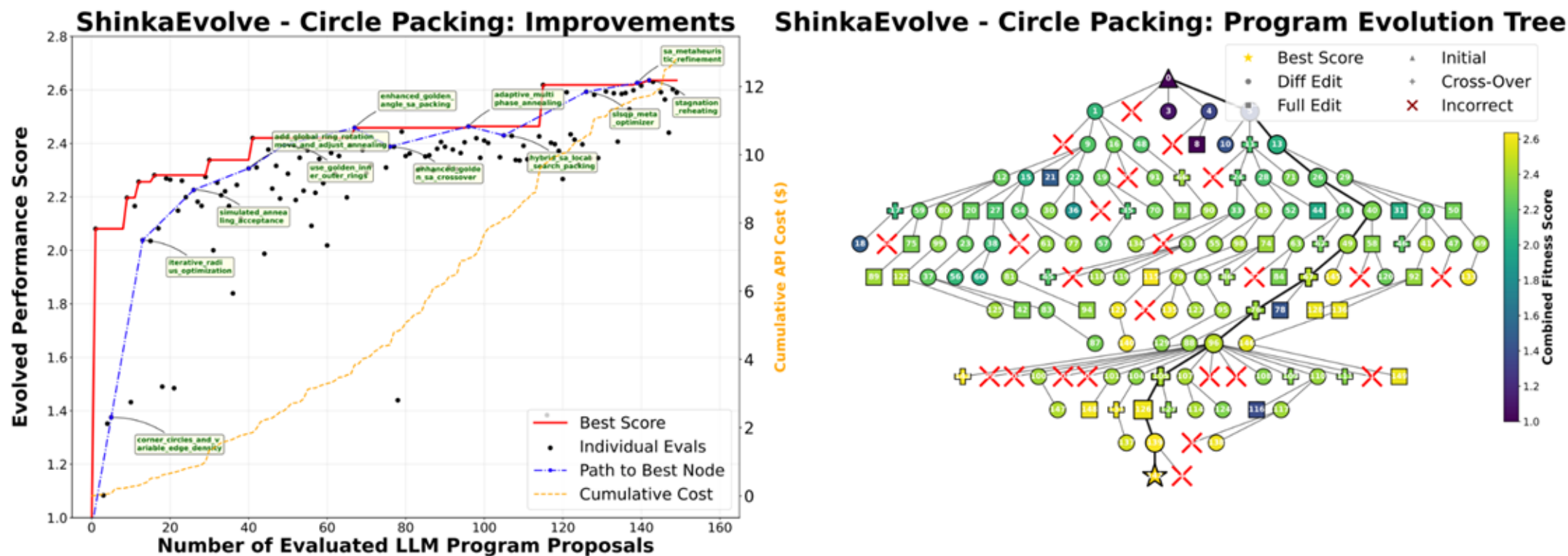
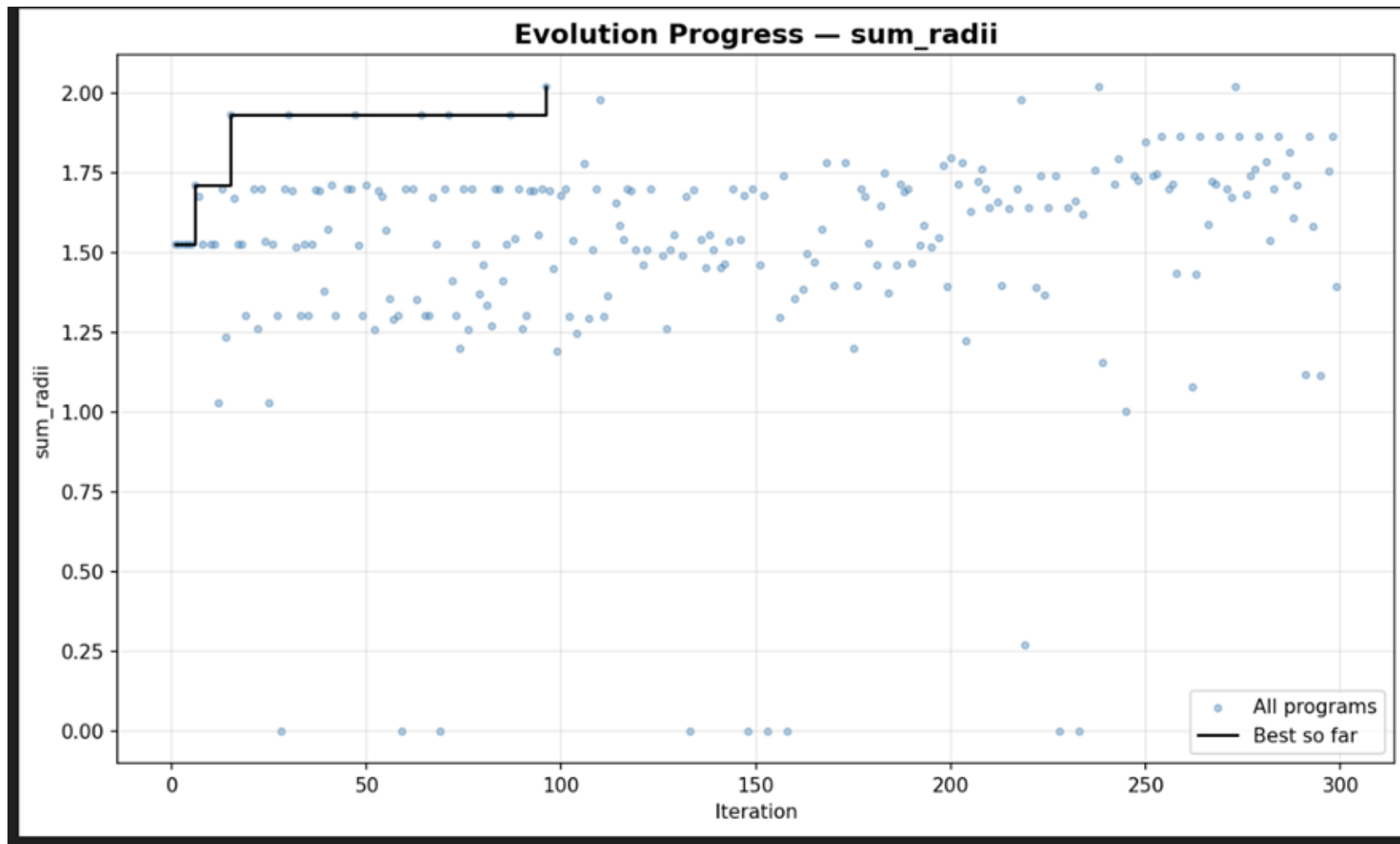


Figure 5 | **SHINKAEVOLVE on Circle Packing Task.** *Left:* SHINKAEVOLVE outperforms AlphaEvolve’s solution within less than 150 program evaluations. *Right:* SHINKAEVOLVE’s program evolution tree demonstrates the iterative composition of stepping stones into high-performing solutions.

Few Samples Lead to Improvements



Evolving the Model

Recent works such as TTT-Discover and ThetaEvolve train the model at test-time with improvement score as the reward signal

Algorithm 1 Test-Time Training to Discover (TTT-Discover)

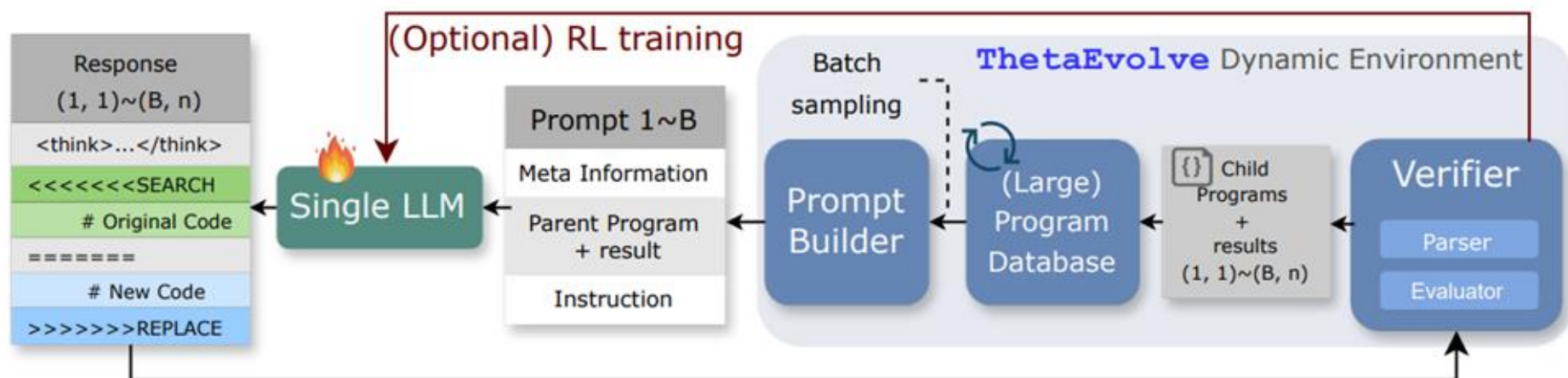
```

1: Input: problem description  $d$  and policy  $\pi_{\theta_0}$  with initial weights  $\theta_0$ .
2:  $R, T = \text{get\_env}(d)$     $\triangleright$   $d$  induces the reward and transition functions of the environment (§2.1)
3:  $\mathcal{H}_0 = \{(\langle \text{empty} \rangle, R(\langle \text{empty} \rangle), \{\})\}$     $\triangleright$  Initialize buffer with the empty solution (§2.2)
4: for  $i = 0, 1, \dots, N - 1$  do
5:    $s_i, c_i \sim \text{reuse}(\mathcal{H}_i)$     $\triangleright$  Sample initial state and context with a reuse heuristic
6:    $a_i \sim \pi_{\theta_i}(\cdot \mid d, s_i, c_i)$     $\triangleright$  Sample action from policy
7:    $s'_i = T(a_i)$     $\triangleright$  Transition to next state
8:    $r_i = R(s'_i)$     $\triangleright$  Evaluate reward of next state
9:    $\mathcal{H}_{i+1} = \mathcal{H}_i \cup \{(s_i, a_i, s'_i, r_i)\}$     $\triangleright$  Add current attempt to buffer
10:   $\theta_{i+1} = \text{train}(\theta_i, (d, s_i, c_i, a_i, r_i))$     $\triangleright$  Improve the model weights with train
11: end for
12: return  $s_{i^*}$ , where  $i^* = \arg \max_{i=0,1,\dots,N-1} r_i$     $\triangleright$  Return the state with the highest reward

```

Evolving the Model

Recent works such as TTT-Discover and ThetaEvolve train the model at test-time with improvement score as the reward signal

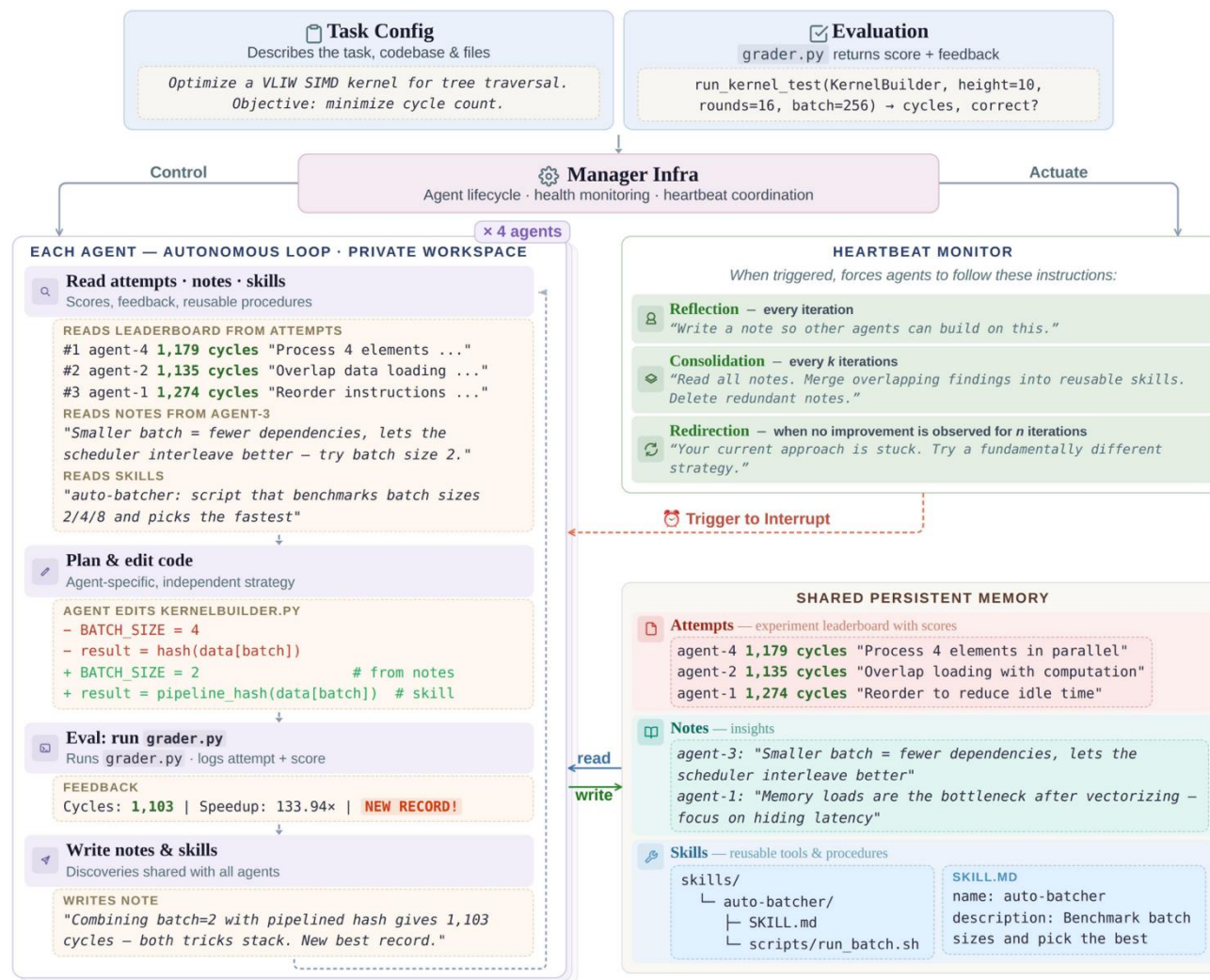


CORAL: Multi-Agent Evolution

<https://github.com/Human-Agent-Society/CORAL>

Key ideas:

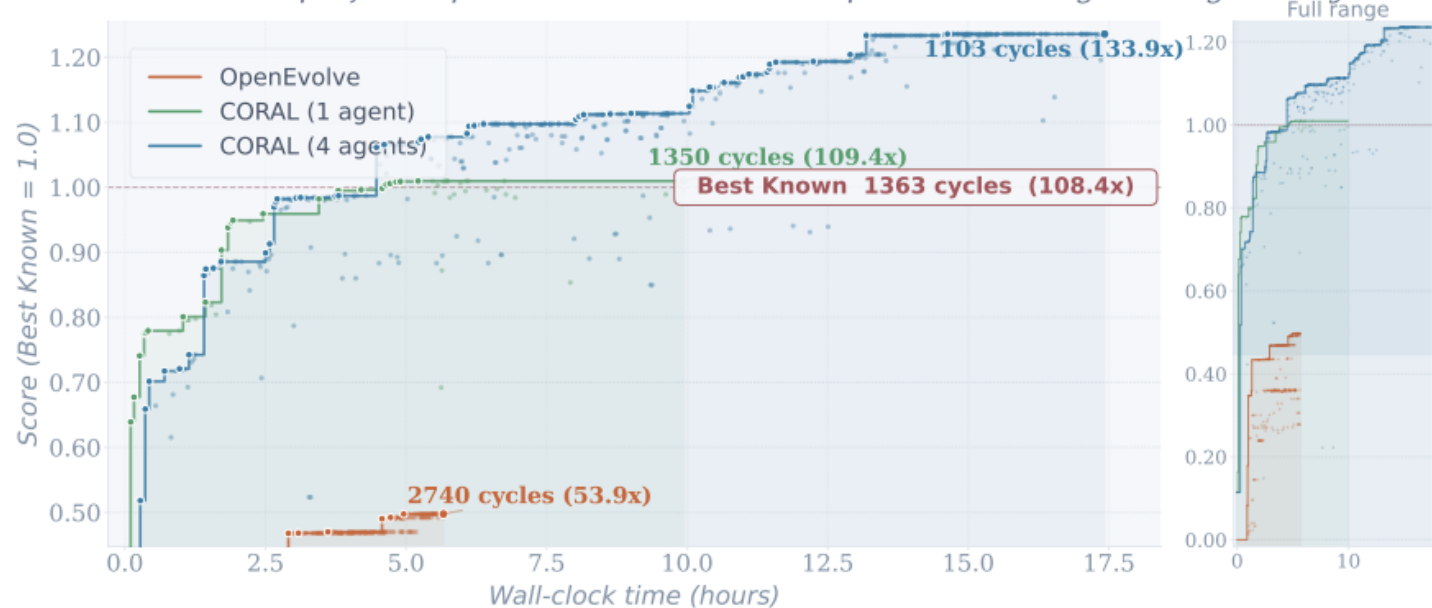
- Multi-agent specialization
- Increased autonomy
- Shared memory
- Heartbeat mechanism



CORAL: Multi-Agent Evolution

<https://github.com/Human-Agent-Society/CORAL>

CORAL outperforms previous SOTA on Anthropic's kernel engineering task by 20%



Operations research

Systems engineering

Scientific discovery

	Model	Evals ↓	Best Cycles ↓	Speedup ↑	Imp. Rate ↑
• Best Known	Anthropic	–	1363	108.4x	–
• OpenEvolve	Opus 4.6	363	2740	53.9x	3% (12/363)
• CORAL (1 agent)	Opus 4.6	56	1350	109.4x	43% (24/56)
• CORAL (4 agents)	Opus 4.6	596	1103	133.9x	9% (54/596)

The task is to optimize a GPU kernel to minimize execution cycles from a 147,734-cycle baseline. Score = baseline / cycles, normalized so the previous best known result (1,363 cycles) equals 1.0. All methods use Claude Opus 4.6.

Self-Evolving Agents

How agents can constantly verify their knowledge and improve their capabilities over time beyond static human supervision.

See blog post: <https://quao627.github.io/blog/self-evolving-agents/>

Assignments for This Coming Week

Please fill out course evaluations and give us feedback!

HW5 due this Friday May 8.

For project:

- Make sure to meet with myself and TAs this week if you need advice.
- Presentations next Tuesday May 12.
- Final report due Tuesday May 19.