# From Multimodal LLM to Human-level AI
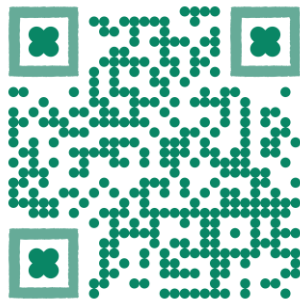
*Architecture*, *Modality*, *Function*, *Instruction*, *Hallucination*, Evaluation, *Reasoning* and **Beyond**
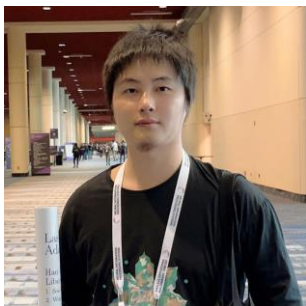
https://mllm2024.github.io/ACM-MM2024/
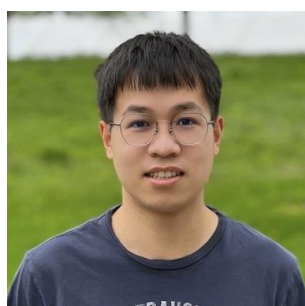
ACM Multimedia 2024

Melbourne, Australia

**Hao Fei**
*National University of Singapore*

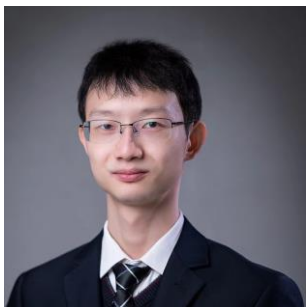**Xiangtai Li**
*ByteDance/Tiktok*

**Haotian Liu**
*xAI*

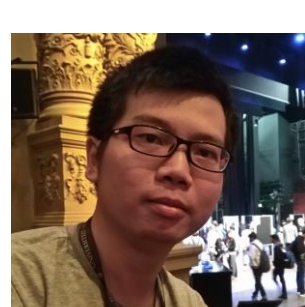**Fuxiao Liu**
*University of Maryland, College Park*

**Zhuosheng Zhang**
*Shanghai Jiao Tong University*

**Hanwang Zhang**
*Nanyang Technological University*

**Kaipeng Zhang**
*Shanghai AI Lab*

**Shuicheng Yan**
*Kunlun 2050 Research, Skywork AI*

# Part-VI
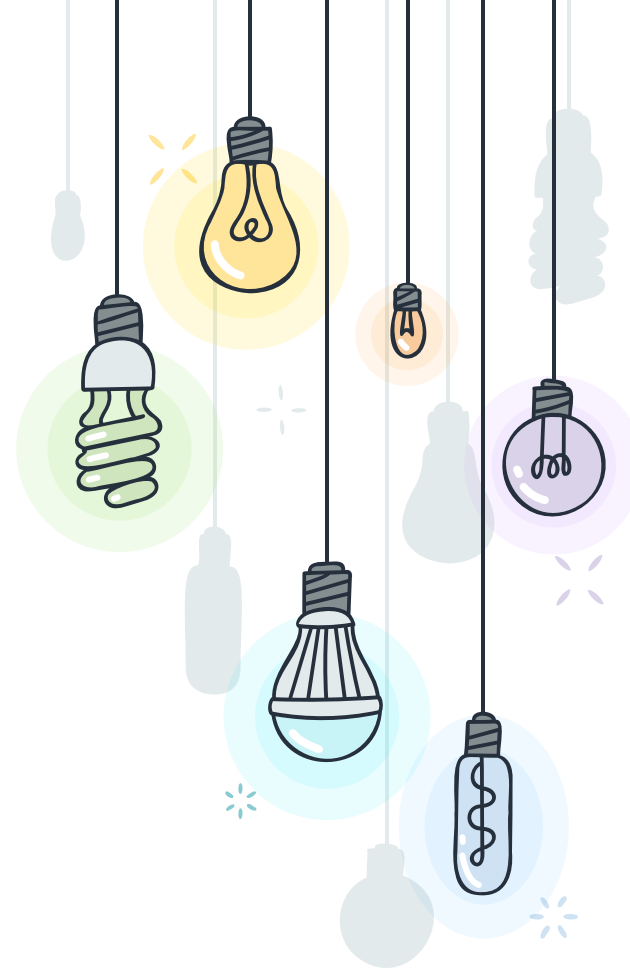
# Evaluation & Generalist
## *Road to L5 MM Generalist*

**Hanwang Zhang**
**Assoc. Prof**

*Nanyang Technological University*

*hanwangzhang@ntu.edu.sg*

https://mreallab.github.io/
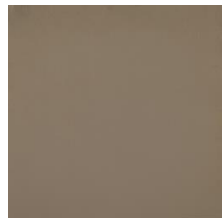
# Road to L₅ MM Generalist

*https://path2generalist.github.io*
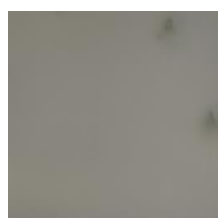
Hanwang Zhang 张含望

Hao Fei
NUS

Yuan Zhou
NTU

Juncheng Li
ZJU

Kaihang Pan
ZJU

Zhongqi Yue
NTU

Shuicheng Yan
Skywork AI

# Comprehension & Generation

(1) Comprehension is a "**many-to-one**" paradigm | (2) Generation is an "**One-to-one**" paradigm

**Example:**

What is it?



Comprehension

Dog

**Example:**

Generate a dog according to descriptions!

description1    description2    description3

Generation



*Inputs and outputs are matched one-by-one.

**Differences between Comprehension and Generation**

# Today's benchmarks are challenging but still fail to systematically reflect MLLMs' <span style="color:red">synergy</span> in/across comprehension and generation.



**MMT-Bench Benchmark**
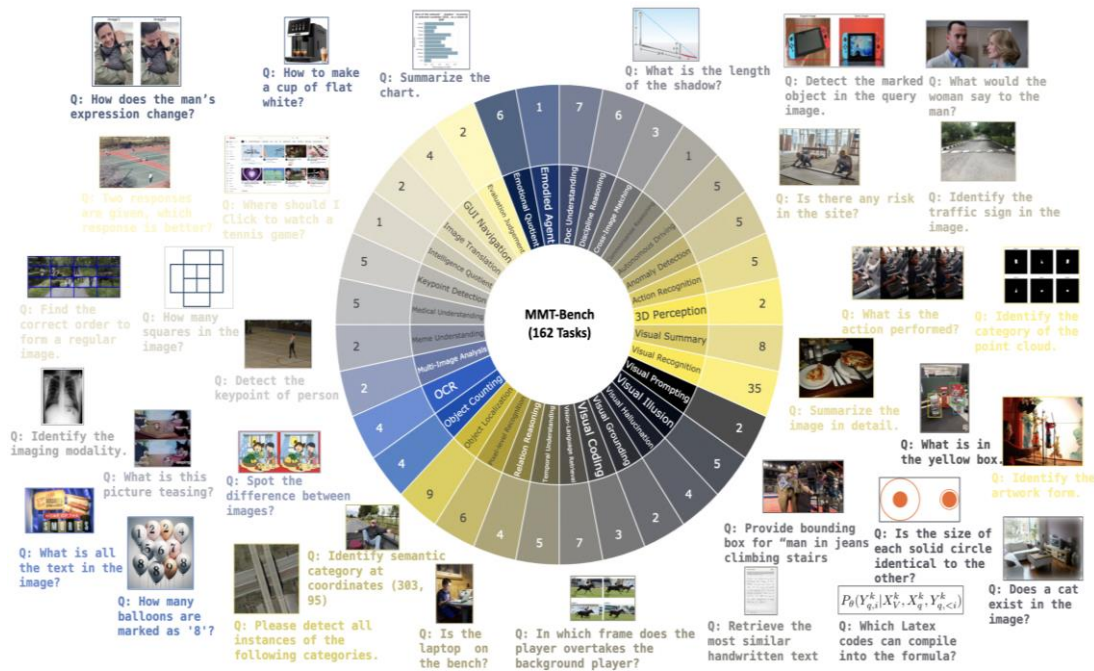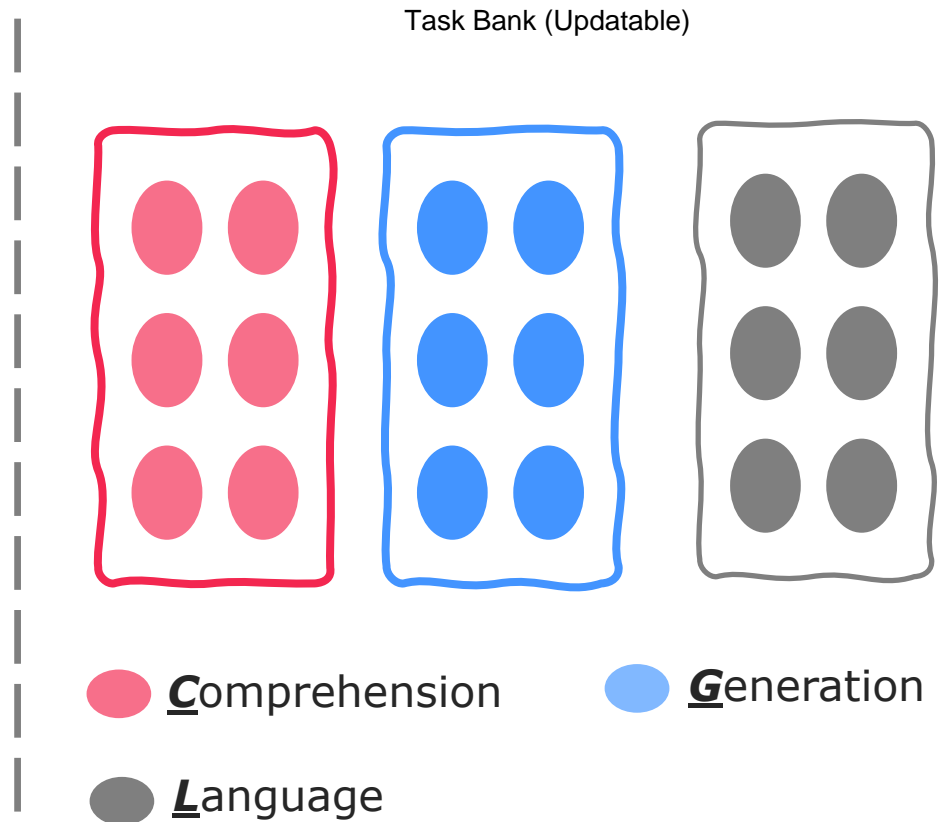
| Rank | Model | Score |
|------|-------|-------|
| 1 | GPT4o | 65.5 |
| 2 | InternVL-Chat-v1.2-34B | 63.4 |
| 3 | QwenVLMax | 62.4 |
| 4 | Qwen-VL-Plus | 62.3 |
| 5 | GeminiProVision | 61.6 |
| 6 | GPT4V_20240409 | 61.1 |
| 7 | LLaVA-NEXT-34B | 60.8 |
| 8 | XComposer2 | 55.7 |
| 9 | BLIP2 | 54.8 |

Level 5
Models are task-unified players, and synergy is across C, G, and L

Level 4
Models are task-unified players, and synergy is across C and G

Level 3
Models are task-unified players, and synergy is in C and G

Level 2
Models are task-unified players

Level 1 (not scored)
Models are task-specific players

Task Bank (Updatable)

**C**omprehension   **G**eneration

**L**anguage

# Task dispatcher is NOT synergy

ViperGPT: Visual Inference via Python Execution for Reasoning. ICCV'23
Visual Programming: Compositional visual reasoning without training. CVPR'23

# Comprehension & Generation

**Comprehension**

FSC147 · PASCAL VOC 2012

CARPK · NYUv2 · OKVQA · MVTecAD

TextOCR · AM-2K · DocVQA · Cityscapes

CIFAR-100 · OCHuman · RefCOCO+ · GQA

Flickr30k · MS COCO · VQAv2 · ADE20K

COCO-caption · RefCOCO · VizWiz

RefCOCOg

**Generation**

MM CelebA-HQ · Places2

Imagenet · COCO-Stuff · Manga109

MSCOCO · Visual Genome · Set14

Scribble · BSD100 · urban · SketchyCOCO

CelebA-HQ · FFHQ · HIDE · LOL · ADE20K

PIE-Bench · VITON-HD · Cityscapes
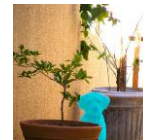
SIDD · GoPro · CUB

**Example:**



# Are there elephants in the image?
# Yes



# Is the answer to the above question 65?
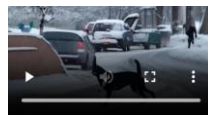# Yes



# Segment out the dog from the image.



# Describe the photo.
# A child in a purple outfit is seated on a chair.

**Comprehension**

**Generation**

**Example:**

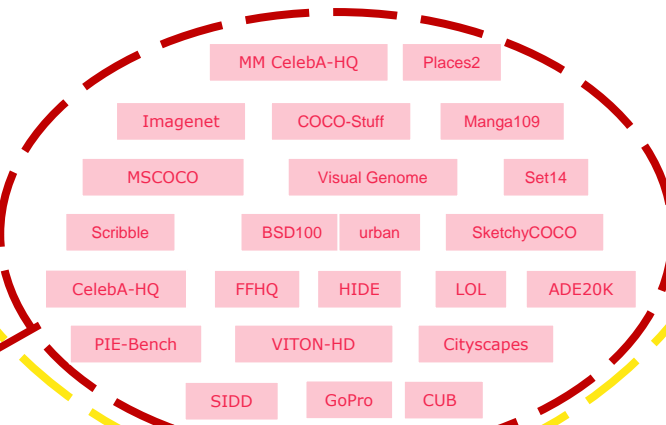# Please generate a video where a dog run past a car on the street in the snow.



# Swap out the background of the video for a snowy scene.





# Generate a picture about burning fire.

# NLP



## Language

| 20 Newgroups | AG News | IMDB | SST | Yelp | TREC | RACE | MultiRC |

| DBpedia | FakeNewsNet | SNLI | Quora | NER | CNN-Daily Mail |

| CoNLL2003 | OntoNotes5.0 | Semeval2010-task-8 | DialogRE | ReCoRD |

| SQuAD2.0 | HotpotQA | CoQA | NewQA | SemEval | SNIPS | MS MARCO |

FLAN-T5-XL

Level 5
Models are task-unified players, and synergy is across C, G, and L

Level 4
Models are task-unified players, and synergy is across C and G

Level 3
Models are task-unified players, and synergy is in C and G

Level 2
Models are task-unified players

Level 1
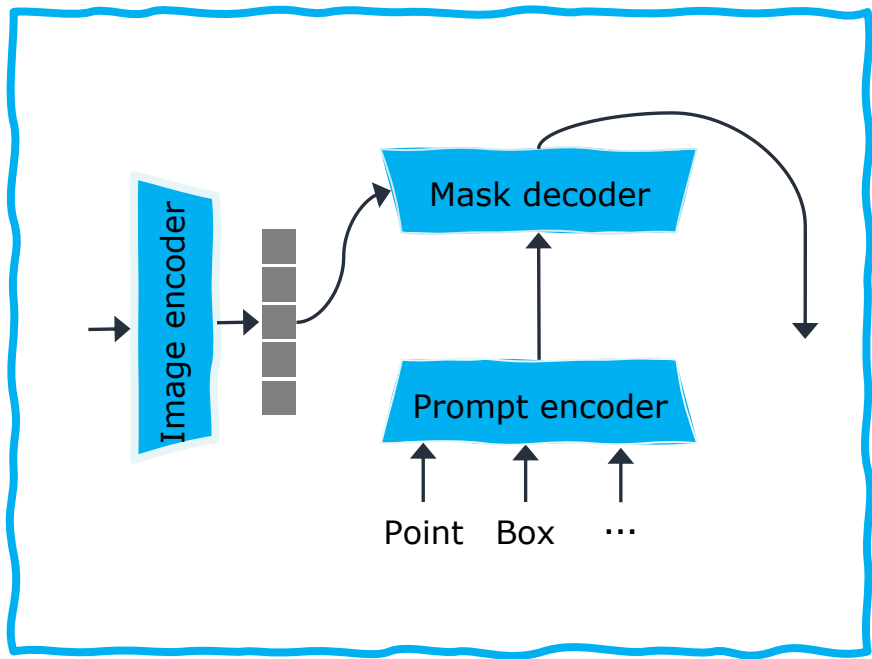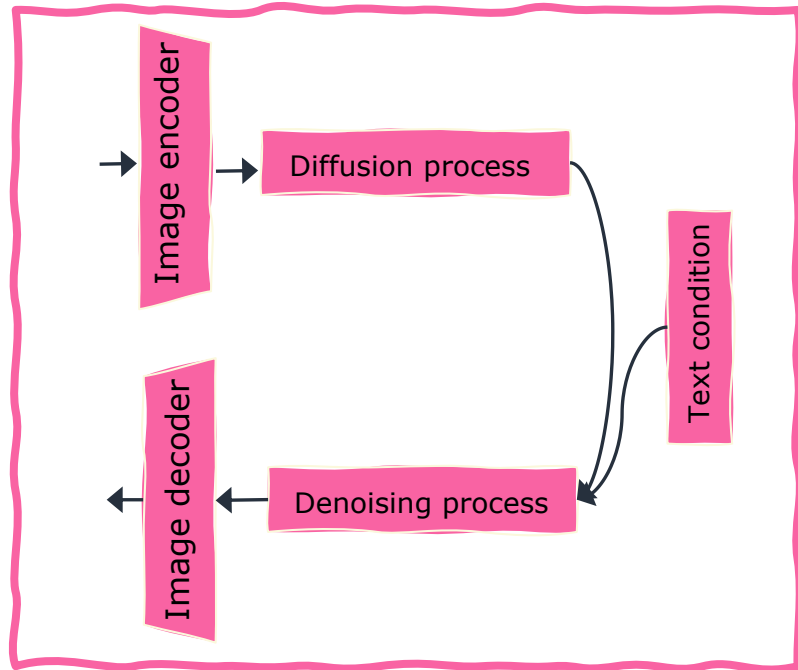Models are task-specific players

Comprehension        Generation

# Level 1: Specialist



SAM (Segmentation)

Stable Diffusion (Text-to-image)

**Examples for the framework of specialist models.
They are specially designed/fine-tuned for specific tasks.**

# Level 1: Specialist

## Comprehension

| FSC147 | PASCAL VOC 2012 | CARPK | NYUv2 | OKVQA |
|---|---|---|---|---|
| counTR | SegCLIP | counTR | TransDepth | GIT |

| TextOCR | AM-2K | DocVQA | Cityscapes | CIFAR-100 |
|---|---|---|---|---|
| Parseq | GFM | Donut | OneFormer | Astroformer |

| RefCOCO+ | GQA | Flickr30k | MS COCO | VQAv2 | ADE20K |
|---|---|---|---|---|---|
| polygon-former | GIT | CLIP | DINO | GIT | OneFormer |

| COCO-caption | RefCOCO | VizWiz | MVTecAD | OCHuman |
|---|---|---|---|---|
| GRIT | polygon-former | GIT | CPR | BUCTD |

## Generation

| MM CelebA-HQ | Places2 | Imagenet | COCO-Stuff | Manga109 |
|---|---|---|---|---|
| Lafite | Mat | VQ-diff | layoutdiff | Hat |

| MSCOCO | Visual Genome | Set14 | SketchyCOCO | ADE20K |
|---|---|---|---|---|
| Lafite | layoutdiff | Hat | Piti | Inade |

| Scribble | BSD100 | urban | Cityscapes | HIDE | LOL | FFHQ |
|---|---|---|---|---|---|---|
| Piti | Hat | Hat | Inade | Lakdnet | WaveNet | Mat |

| CelebA-HQ | PIE-Bench | VITON-HD | SIDD | GoPro | CUB |
|---|---|---|---|---|---|
| Mat | P2P | Mgd | Hinet | Lakdnet | Lafite |

**Level 5**
Models are task-unified players, and synergy is across C, G, and L

**Level 4**
Models are task-unified players, and synergy is across C and G

**Level 3**
Models are task-unified players, and synergy is in C and G

**Level 2**
Models are task-unified players

**Level 1**
Models are task-specific players

Comprehension          Generation

$$S_2 = \frac{1}{M+N} \sum_{i=1}^{M+N} \sigma_i$$

# Level 2: Unified C and G



**Unified comprehension**

Large Language model

Image encoder

Text

Image

Flamingo [58]

(QA: Text+image)

**Unified comprehension and Generation**

Image decoder

Large Language model

Image encoder

Text

Image

MM-Interleaved[68]

(QA: Text+image; Text-to-image, image editing)

**Examples for the framework of unifying C and/or G.**

Level 5
Models are task-unified players, and synergy is across C, G, and L

Level 4
Models are task-unified players, and synergy is across C and G

Level 3
Models are task-unified players, and synergy is in C and G

Level 2
Models are task-unified players

Level 1
Models are task-specific players

Comprehension          Generation

$$S_3 = \frac{1}{M+N} \sum_{i=1}^{M+N} \begin{cases} \sigma_i, \sigma_i \geq \sigma_i^{sota} \\ 0, else \end{cases}$$

# The score of Level 3

Level 5
Models are task-unified players, and synergy is across C, G, and L

Level 4
Models are task-unified players, and synergy is across C and G

Level 3
Models are task-unified players, and synergy is in C and G

Level 2
Models are task-unified players

Level 1
Models are task-specific players

Comprehension    Generation

$$S_4 = \frac{2S_G S_C}{S_G + S_C}, \text{where}$$

$$S_G = \frac{1}{M}\sum_{i=1}^{M}\begin{cases}\sigma_i, \sigma_i \geq \sigma_i^{\text{sota}}\\ 0, \text{else}\end{cases},$$

$$S_C = \frac{1}{N}\sum_{j=1}^{N}\begin{cases}\sigma_j, \sigma_j \geq \sigma_j^{\text{sota}}\\ 0, \text{else}\end{cases}$$

Level 5
Models are task-unified players, and synergy is across C, G, and L

Level 4
Models are task-unified players, and synergy is across C and G

Level 3
Models are task-unified players, and synergy is in C and G

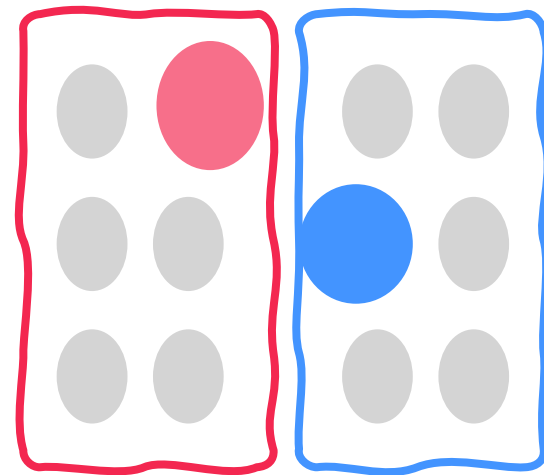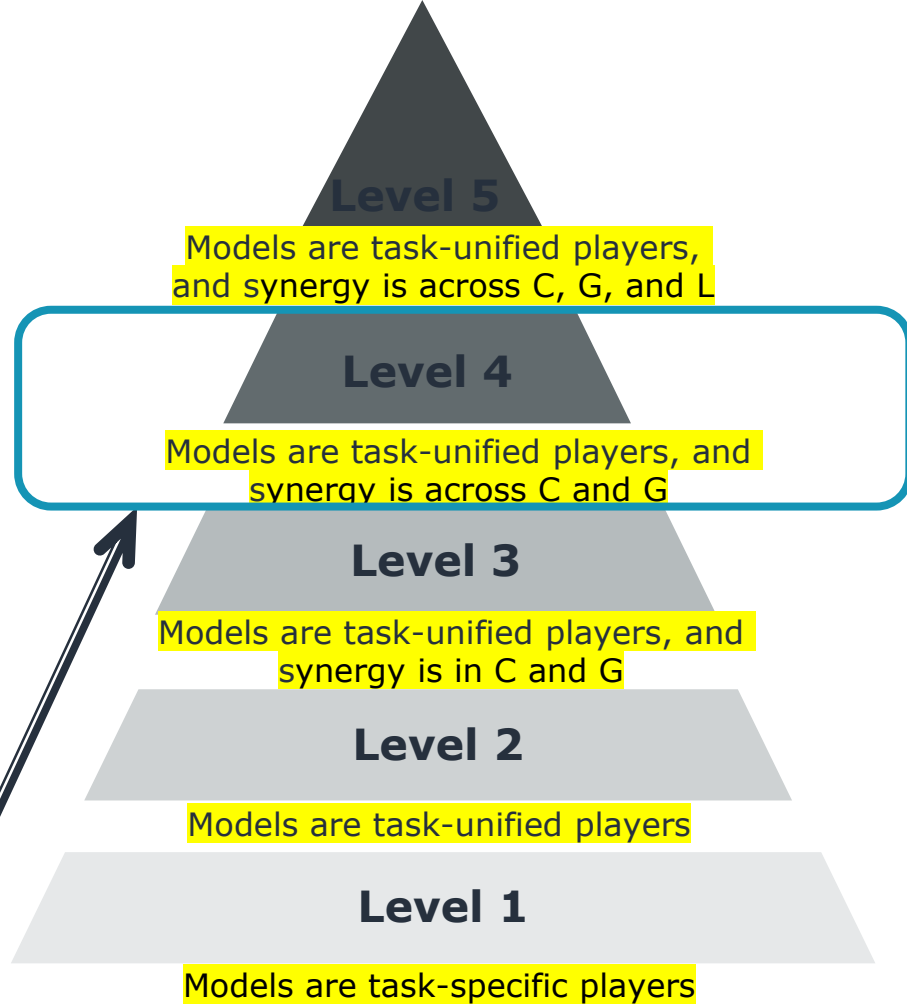Level 2
Models are task-unified players

Level 1
Models are task-specific players

Comprehension    Generation

Language

$$S_5 = S_4 * w_L, \text{where } w_L = \frac{S_L}{100},$$

$$S_L = \frac{1}{T} \sum_{k=1}^{T} \begin{cases} \sigma_k, \sigma_k \geq \sigma_k^{sota} \\ 0, \text{else} \end{cases}$$

**This is our goal!** Level 5: Total Synergy

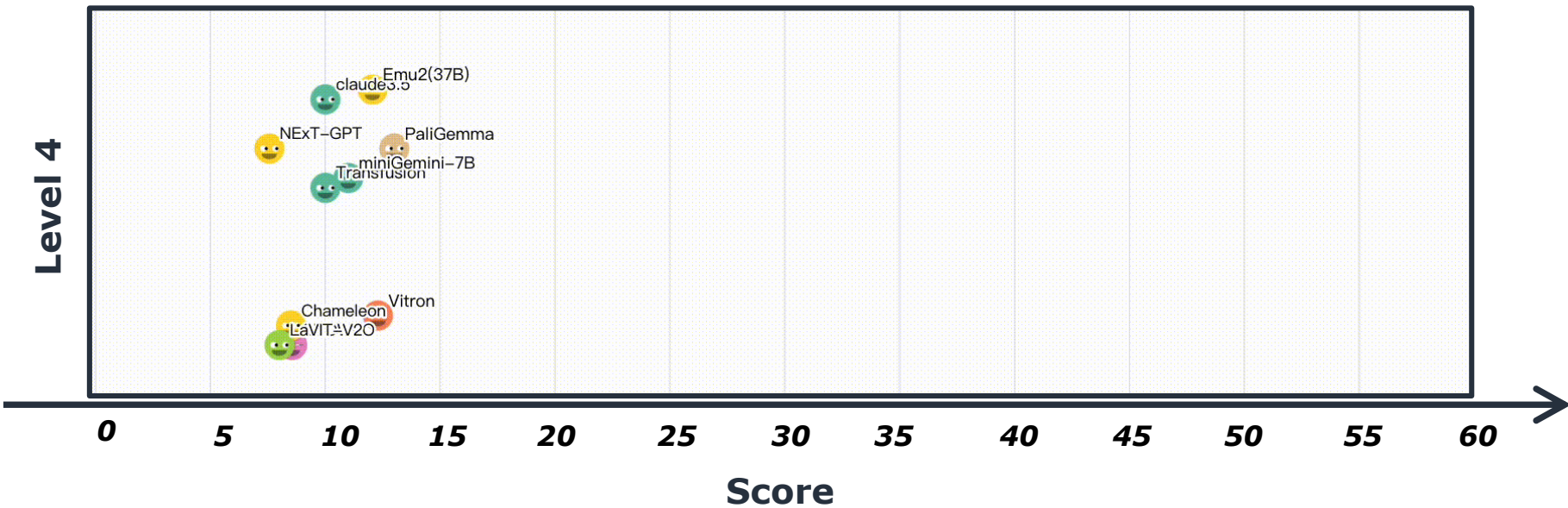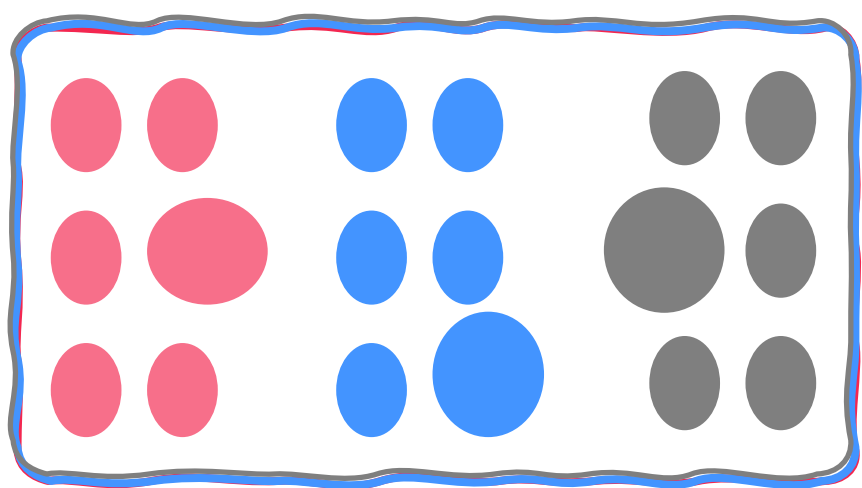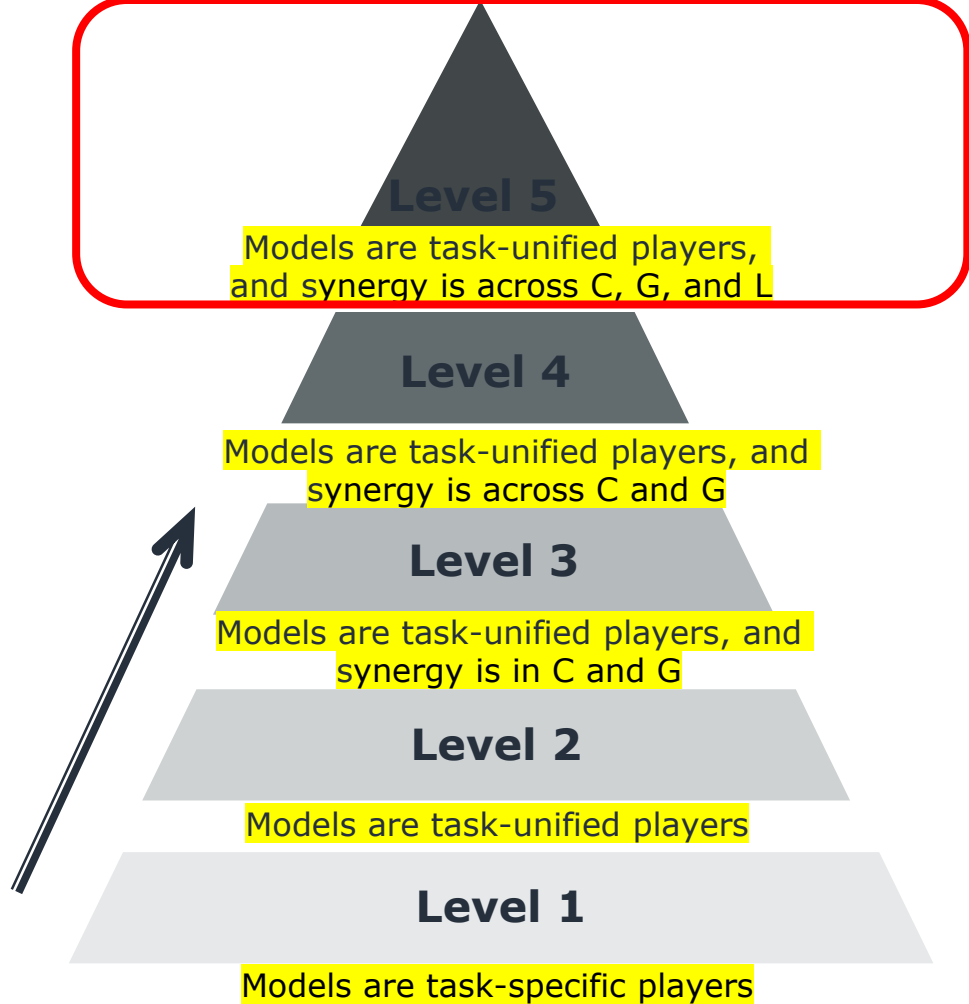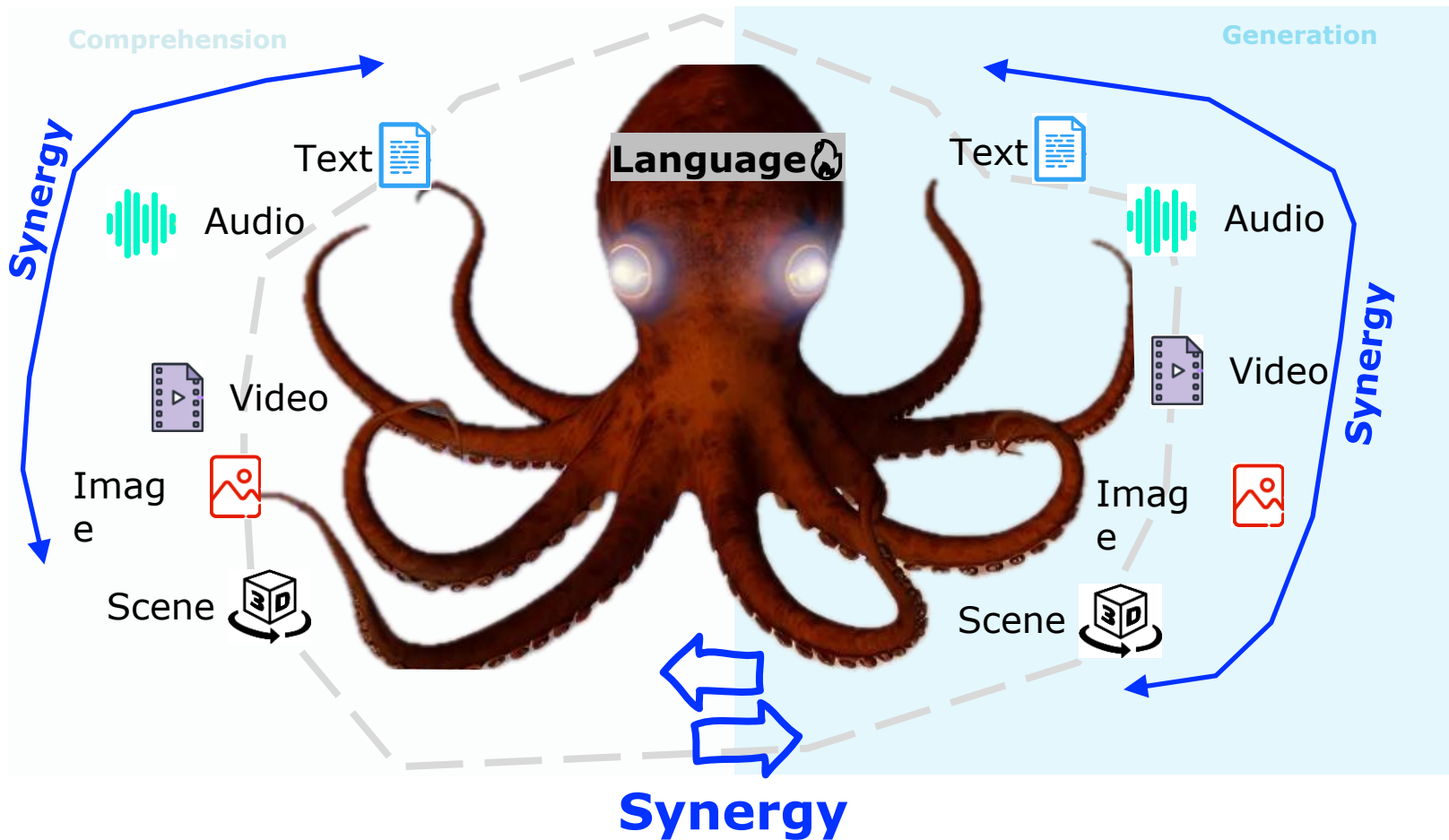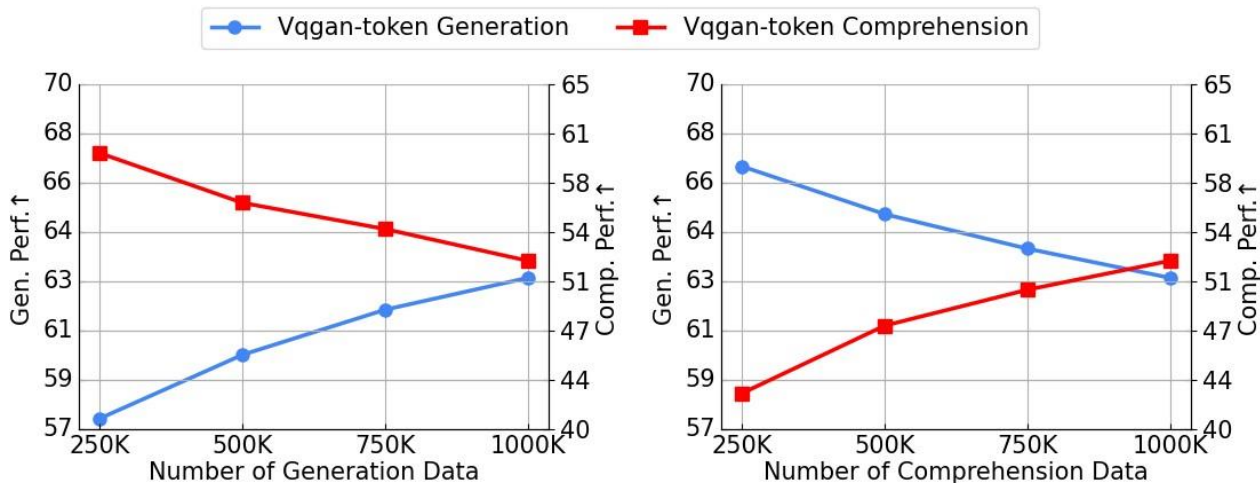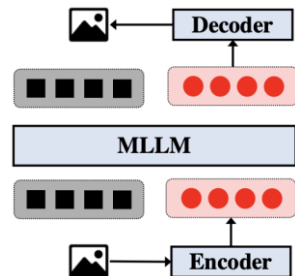| Level | Definition | Score | Example |
|---|---|---|---|
| 1: Specialist | Models are task-specific players | N. A. | **Dino, polygon-former, SegCLIP** |
| 2: Unified C ***and*** G | Models are task-unified players | $S_2 = \dfrac{1}{M+N}\sum_{i=1}^{M+N}\sigma_i$ | **miniGPT4, NextChat** |
| 3: Synergy in C ***and*** G | Models are task-unified players, and synergy is in C and/or G | $S_3 = \dfrac{1}{M+N}\sum_{i=1}^{M+N}\begin{cases}\sigma_i, \sigma_i \geq \sigma_i^{sota}\\0, \text{else}\end{cases}$ | **GPT4v, GPT4o, LLaVA1.5-7b, Qwen-VL-Pius, InternVL, MoE-LLaVA-1.8B-4e, Yi-vl, SEED-LLaMA-14B-SFT, Osprey, GlaMM** |
| 4: Synergy ***across*** C and G | Models are task-unified players, and synergy is across C and G | $S_4 = \dfrac{2S_G S_C}{S_G + S_C}, \text{where}$ <br> $S_G = \dfrac{1}{M}\sum_{i=1}^{M}\begin{cases}\sigma_i, \sigma_i \geq \sigma_i^{sota}\\0, \text{else}\end{cases},$ <br> $S_C = \dfrac{1}{N}\sum_{j=1}^{N}\begin{cases}\sigma_j, \sigma_j \geq \sigma_j^{sota}\\0, \text{else}\end{cases}$ | **miniGemini-7B, Emu2-37B, Vitron,Next-GPT,LaVIT-V2-7B, SHOW-O, Claude3.5 Chameleon, PaliGemma, Transfusion** |
| 5: Total Synergy Synergy ***across*** C, G, and L | Models are task-unified players, and synergy is across C, G, and L | $S_5 = S_4 * w_L, \text{where } w_L = \dfrac{S_L}{100},$ <br> $S_L = \dfrac{1}{T}\sum_{k=1}^{T}\begin{cases}\sigma_k, \sigma_k \geq \sigma_k^{sota}\\0, \text{else}\end{cases}$ | **None, this is our goal!** |

*Upgrade*

**Table1. The roadmap to L5 MM Generalist**

# Climbing: L3->L5

- Most are <=L3

- Why?
  - × Comprehension: V-token to **lose** info to match T-token
  - × Generation: V-token must **preserve** info
  - × V-Language **!=** T-Language
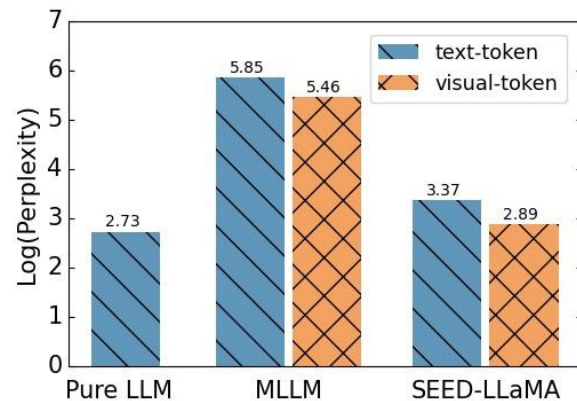  - × V-Generation **!=** T-Generation

# Comprehension and Generation Inconsistency

+ ## Conflicting objectives



Auto-Encoding Morph-Tokens for Multimodal LLM. ICML'24 spotlight

Spatial visual tokens are just word spelling, not language

Noam Chomsky
1928-present

A man

A [*old white*] man

A [old white] man [*with white hair*]

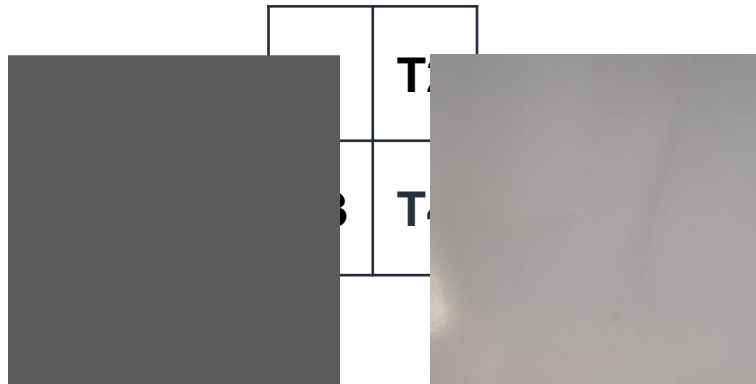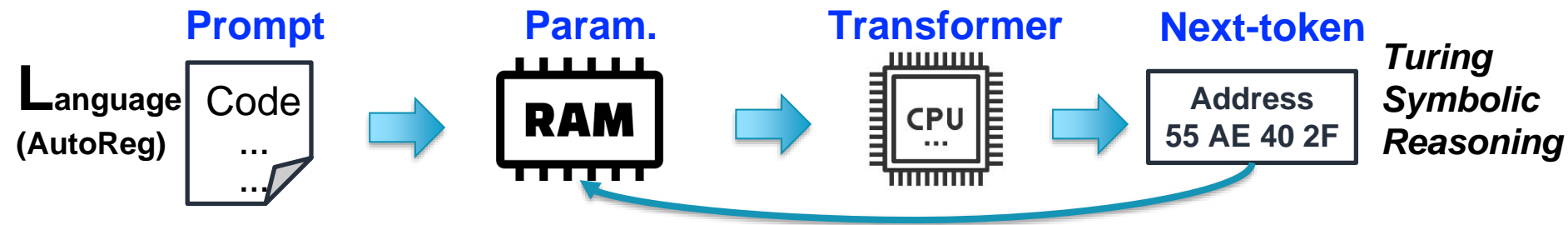A [old white] man [with white hair] [*in black clothes*]

⋮

***Recursive Syntax***
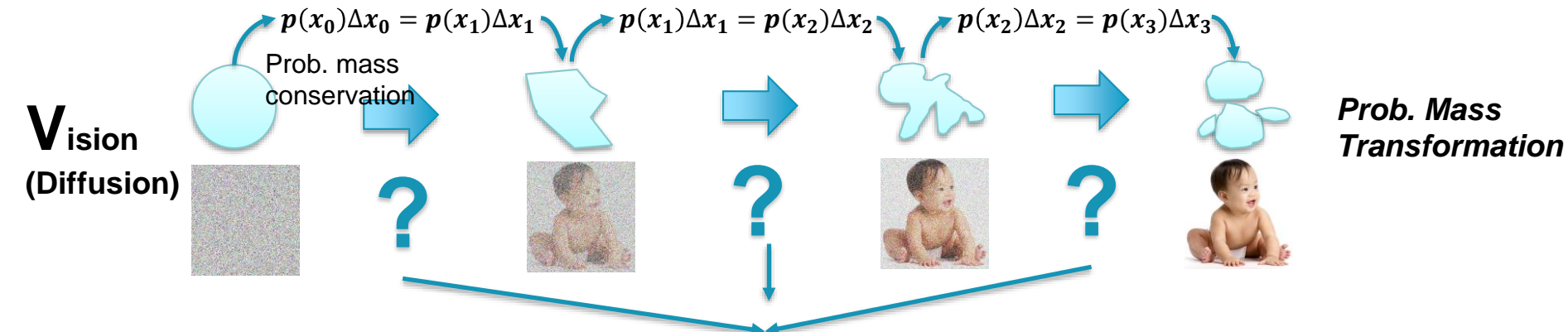
***A furry dog sitting in a striped sofa.***

Whole image is "spelled" as T1 T2 T3 T4

T2

T4

# Challenge 2: Generation (V) ≠ Generation (L)



**Prompt** · **Param.** · **Transformer** · **Next-token**

**L**anguage (AutoReg)

Code ... ...

**RAM**

**CPU** ...

Address 55 AE 40 2F

*Turing Symbolic Reasoning*

≠ *Training objectives never align* ☹

**V**ision (Diffusion)

$p(x_0)\Delta x_0 = p(x_1)\Delta x_1$    $p(x_1)\Delta x_1 = p(x_2)\Delta x_2$    $p(x_2)\Delta x_2 = p(x_3)\Delta x_3$

Prob. mass conservation

*Prob. Mass Transformation*

? ? ?

*the **cause** of transformation?*

# Thank you

Road to L5 MM Generalist

https://path2generalist.github.io