

---

---

# CLiMB



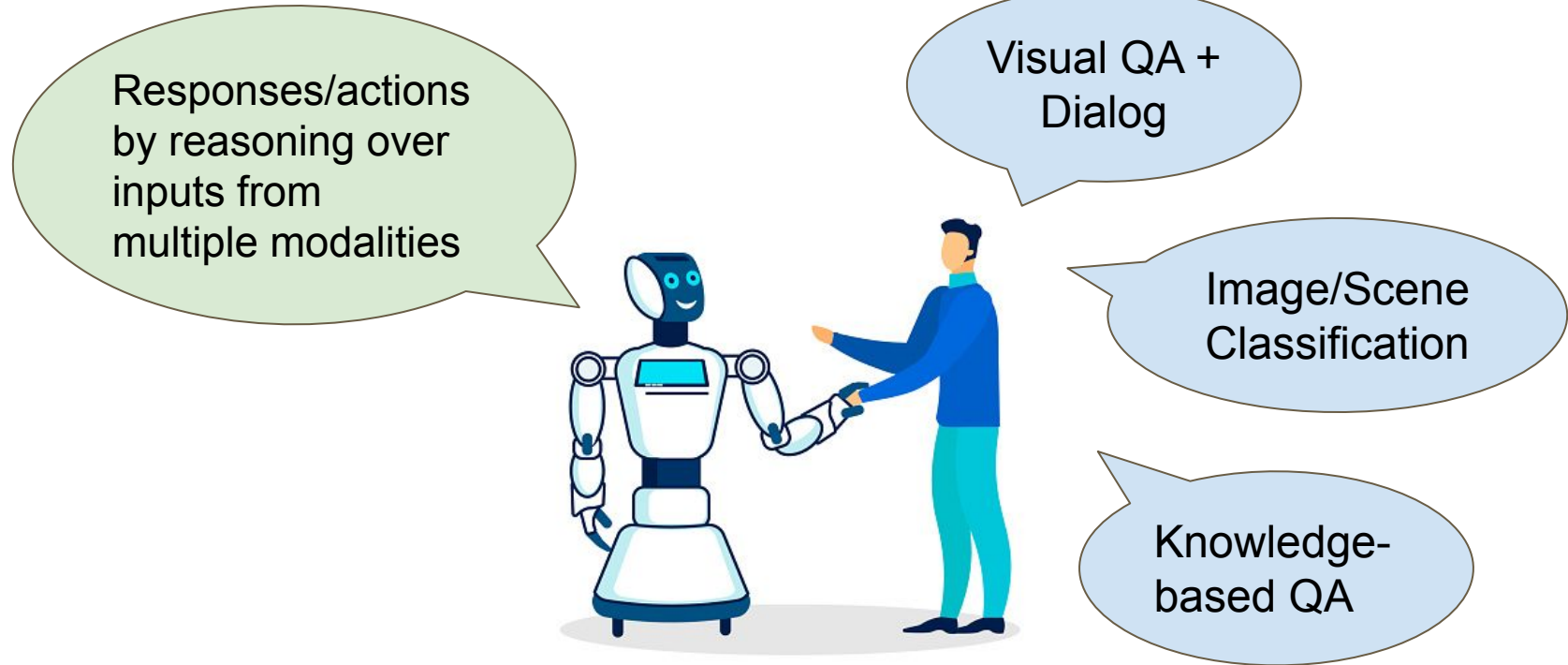
## A Continual Learning Benchmark for Vision-and-Language Tasks

Tejas Srinivasan, Ting-Yun Chang, Leticia Pinto-Alva,  
Georgios Chochlakis, Mohammad Rostami, Jesse Thomason

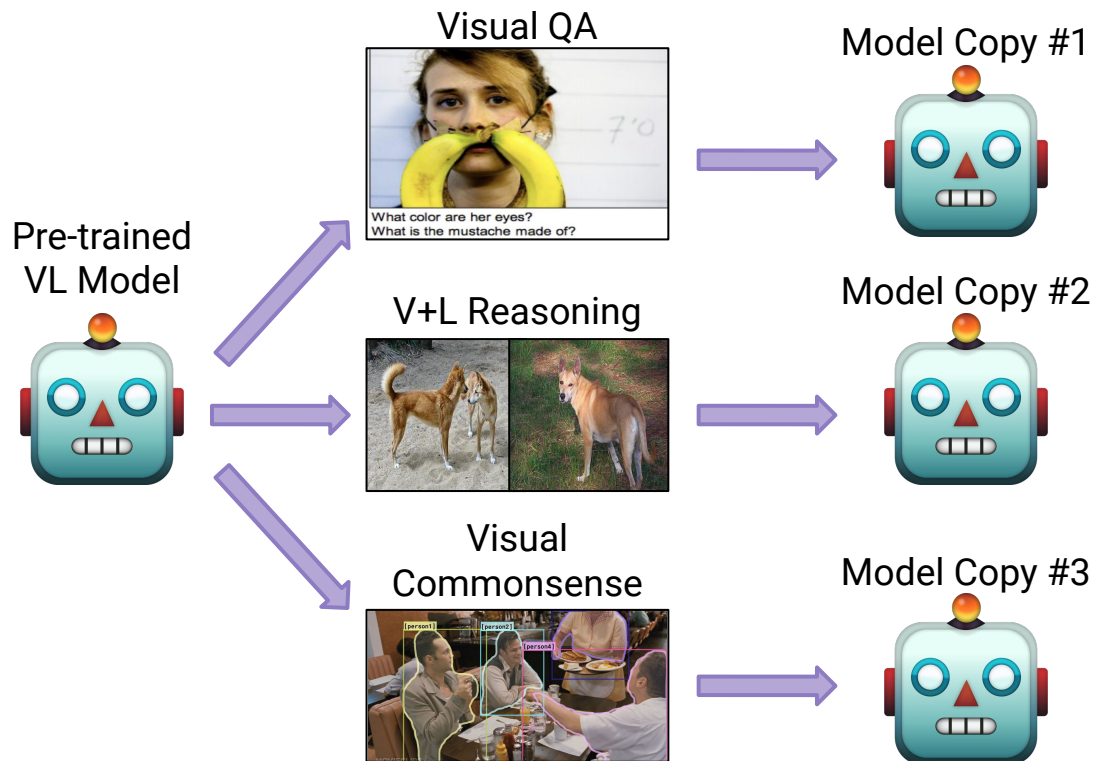
---

---

# Multimodal Agents that can be Deployed

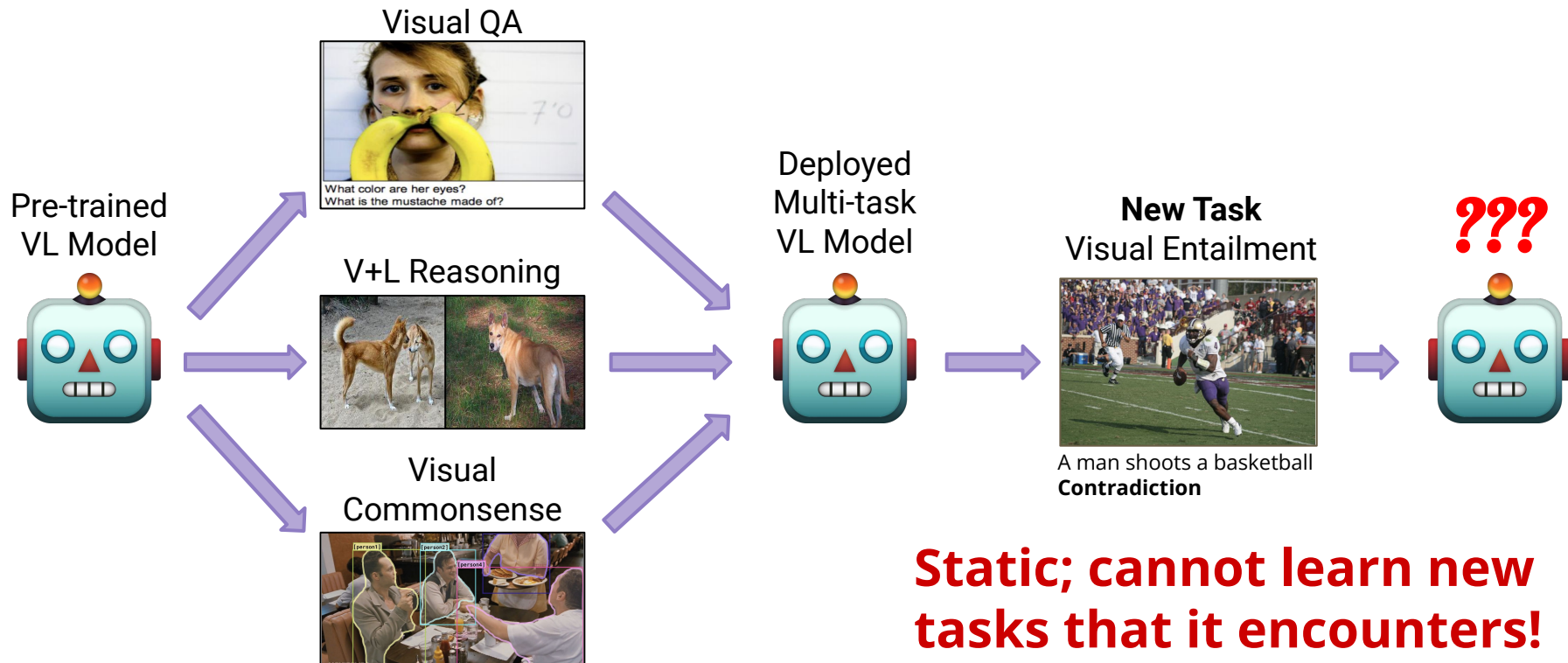


# Paradigms of VL Deployment: Single-Task Finetuning



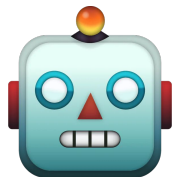
**Need to store a copy of the model for each task!**

# Paradigms of VL Deployment: Multi-Task Learning



# Paradigms of VL Deployment: Continual Learning

Pre-trained  
VL Model



VQA

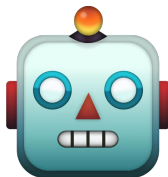


What color are her eyes?  
What is the mustache made of?

NLVR2



...

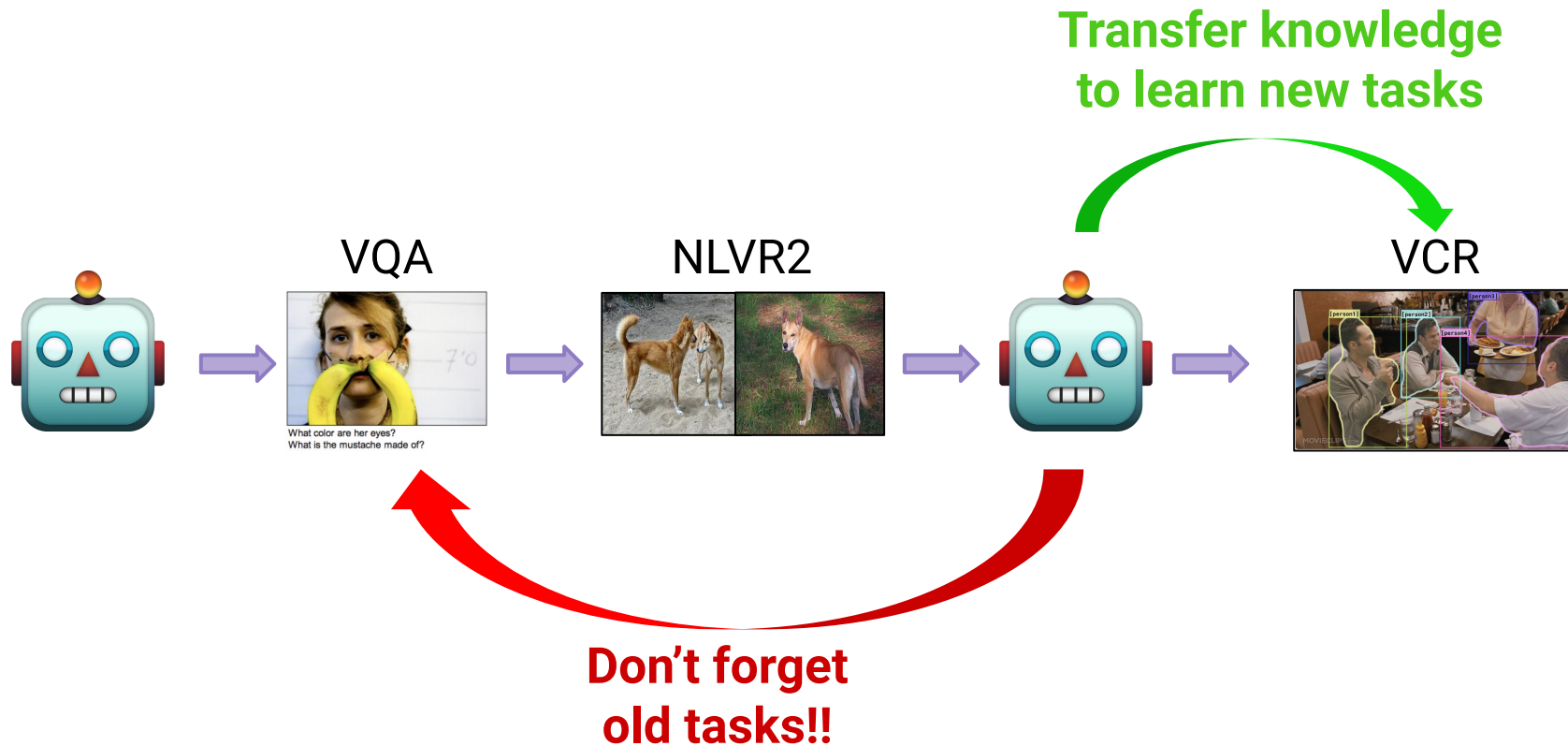


VCR



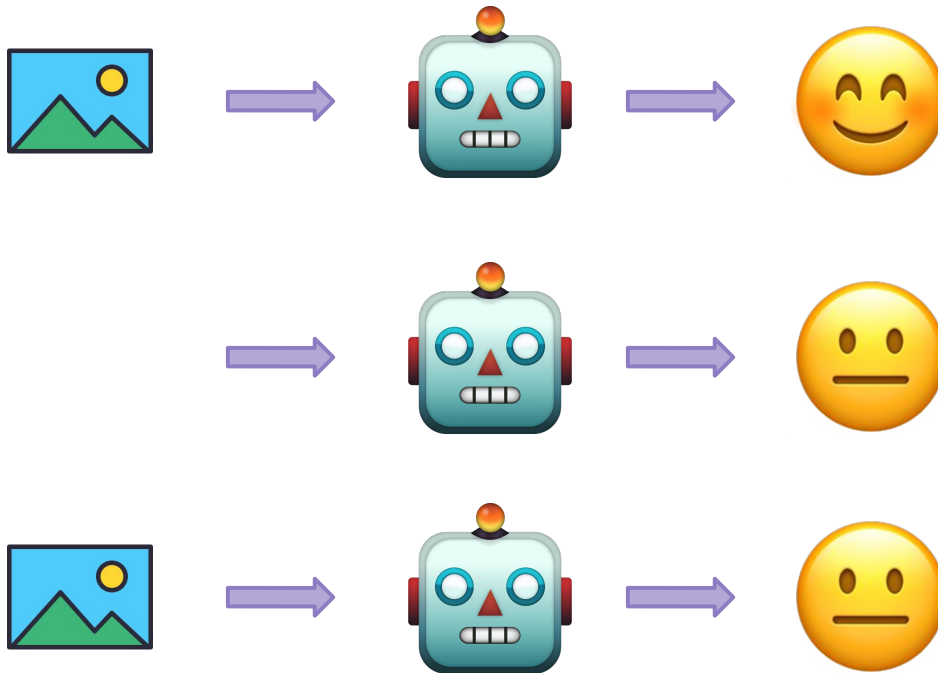
**Dynamic, continually evolving paradigm**  
**Unexplored in multimodal domain!**

# Challenges of Multimodal Continual Learning Deployment

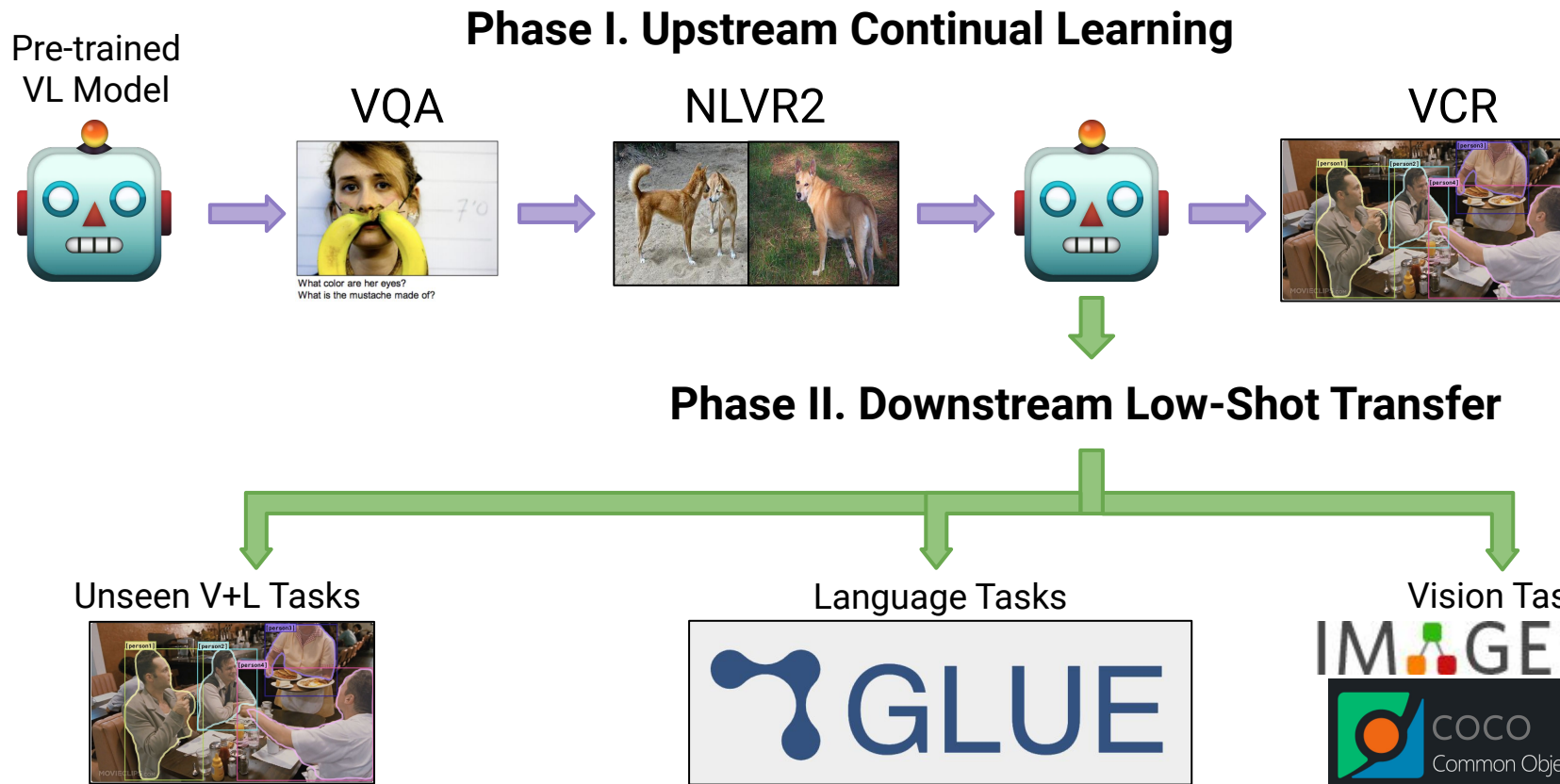


# Challenges of Multimodal Continual Learning Deployment

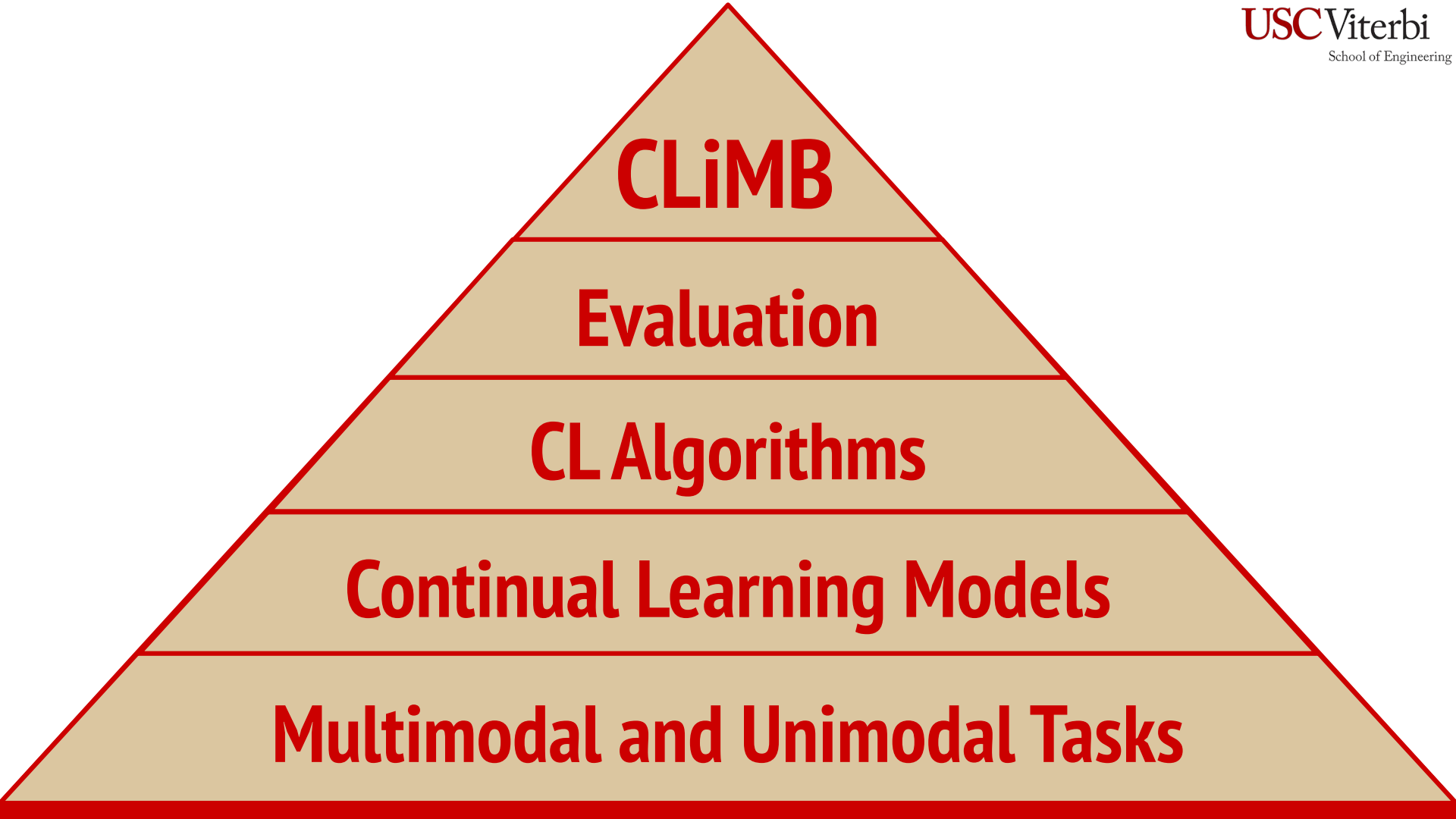
**Not guaranteed to have all modalities when encountering new tasks!**



# CLiMB: The Continual Learning in Multimodality Benchmark





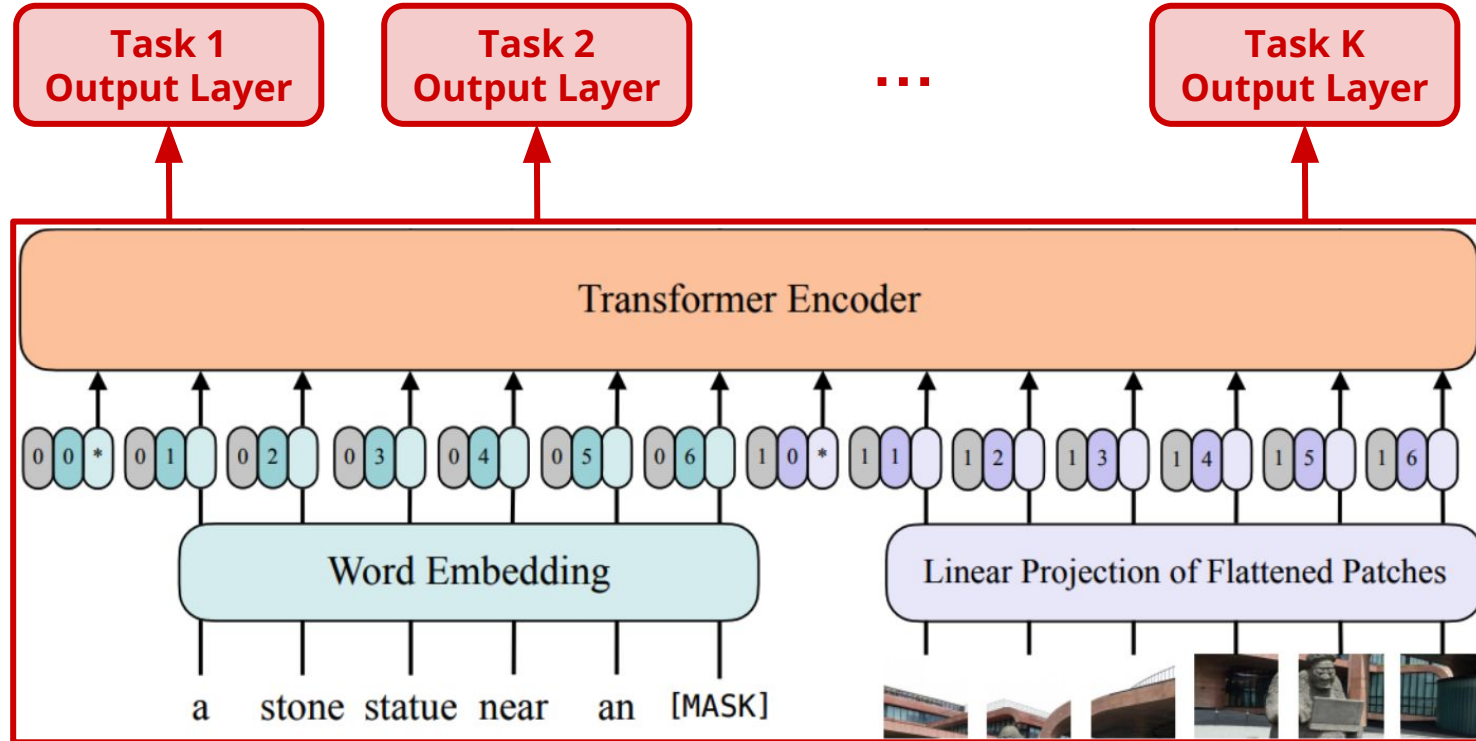


# I. Multimodal and Unimodal Tasks

<b>Vision-and-Language Tasks</b>	<ul style="list-style-type: none"><li>• Visual Question Answering (VQAv2)</li><li>• Natural Language Visual Reasoning (NLVR2)</li><li>• Visual Entailment (SNLI-VE)</li><li>• Visual Commonsense Reasoning (VCR)</li></ul>
<b>Language-Only Tasks</b>	<ul style="list-style-type: none"><li>• IMDb, SST-2 Sentiment Classification</li><li>• HellaSwag</li><li>• CommonsenseQA</li><li>• Physical Interaction QA (PIQA)</li></ul>
<b>Vision-Only Tasks</b>	<ul style="list-style-type: none"><li>• ImageNet-1K Image Classification</li><li>• iNaturalist2019 Image Classification</li><li>• Places365 Image Classification</li><li>• MS-COCO Object Detection</li></ul>

**CLIMB can be easily extended to include new multimodal and unimodal tasks!**

## II. Continual Learning Models



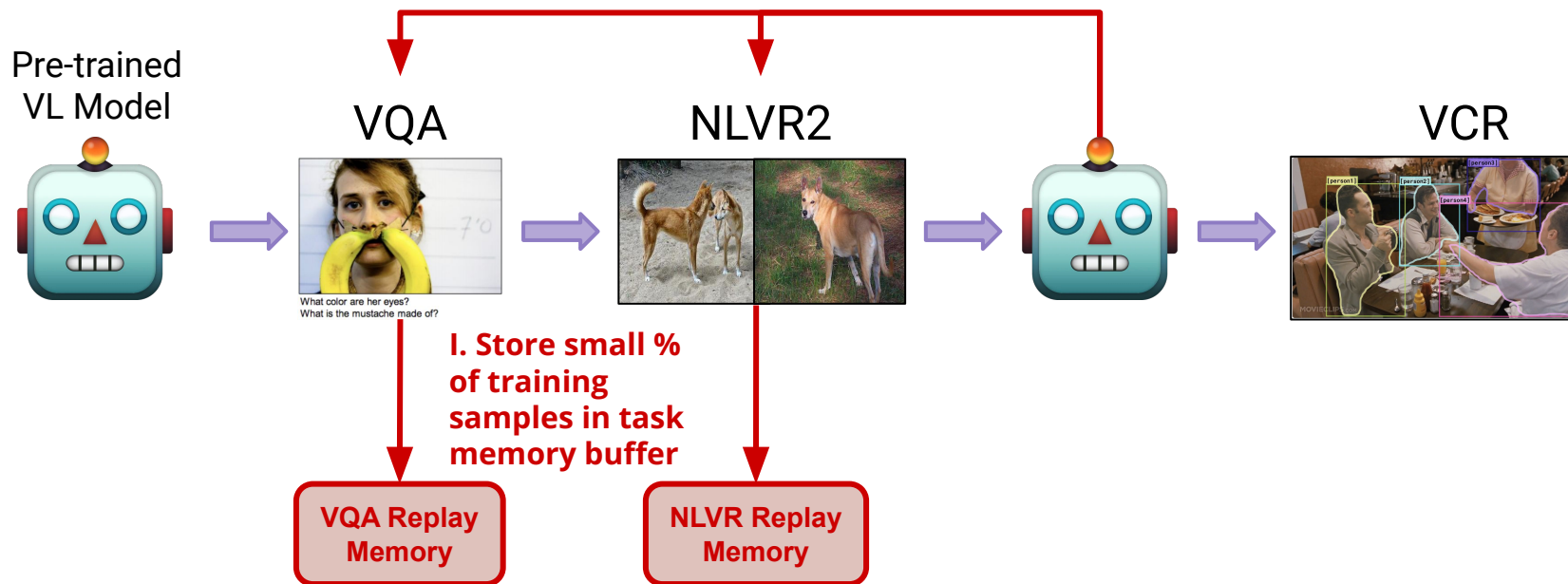
# III. Continual Learning Algorithms

Currently, CLiMB supports 6 different Continual Learning algorithms:

- **Sequential Fine-tuning:** Fine-tune full encoder and task-specific layers
- **Frozen Encoder:** Train only task-specific layers
- **Frozen Bottom-K:** Fine-tune only top encoder layers and task layers
  - We set  $K=9$
- **Experience Replay (ER)**
- **Elastic Weight Consolidation (EWC)**
- **Adapters**

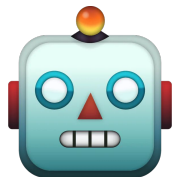
# Experience Replay

**II. Periodically replay a batch from one of the previous task's buffers**



# Elastic Weight Consolidation

Pre-trained  
VL Model



VQA



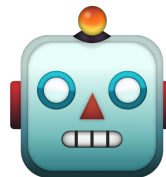
What color are her eyes?  
What is the mustache made of?

NLVR2

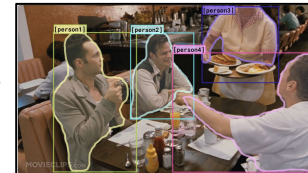


**I. Store previous  
task's model weights**

**NLVR  
Encoder Ckpt**



VCR

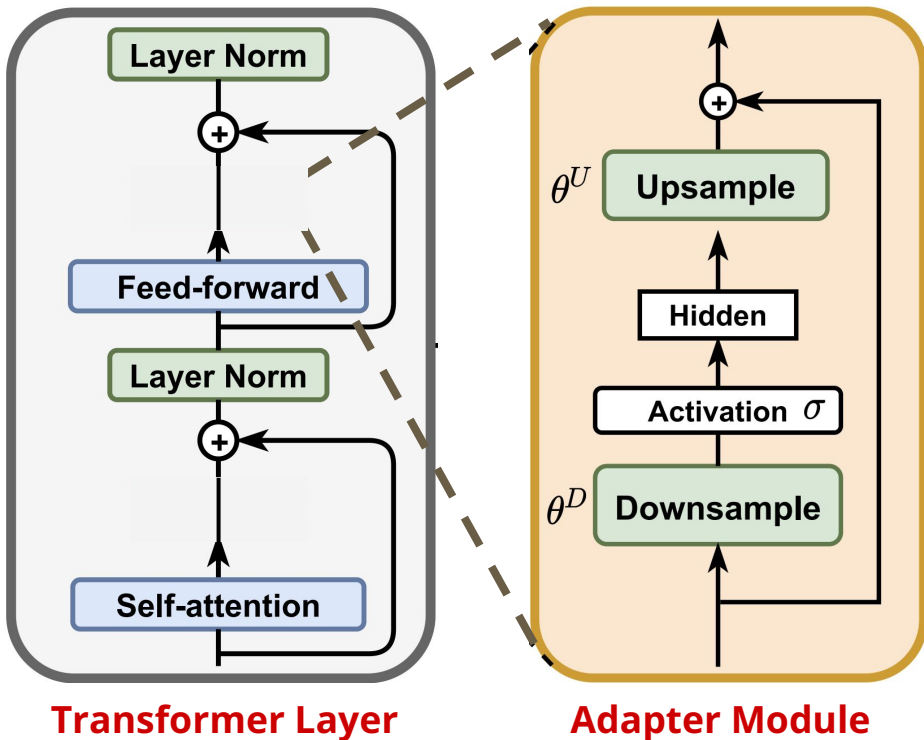


**II. When training on new  
task, add L2 loss between  
model's weights and last  
saved ckpt weights**

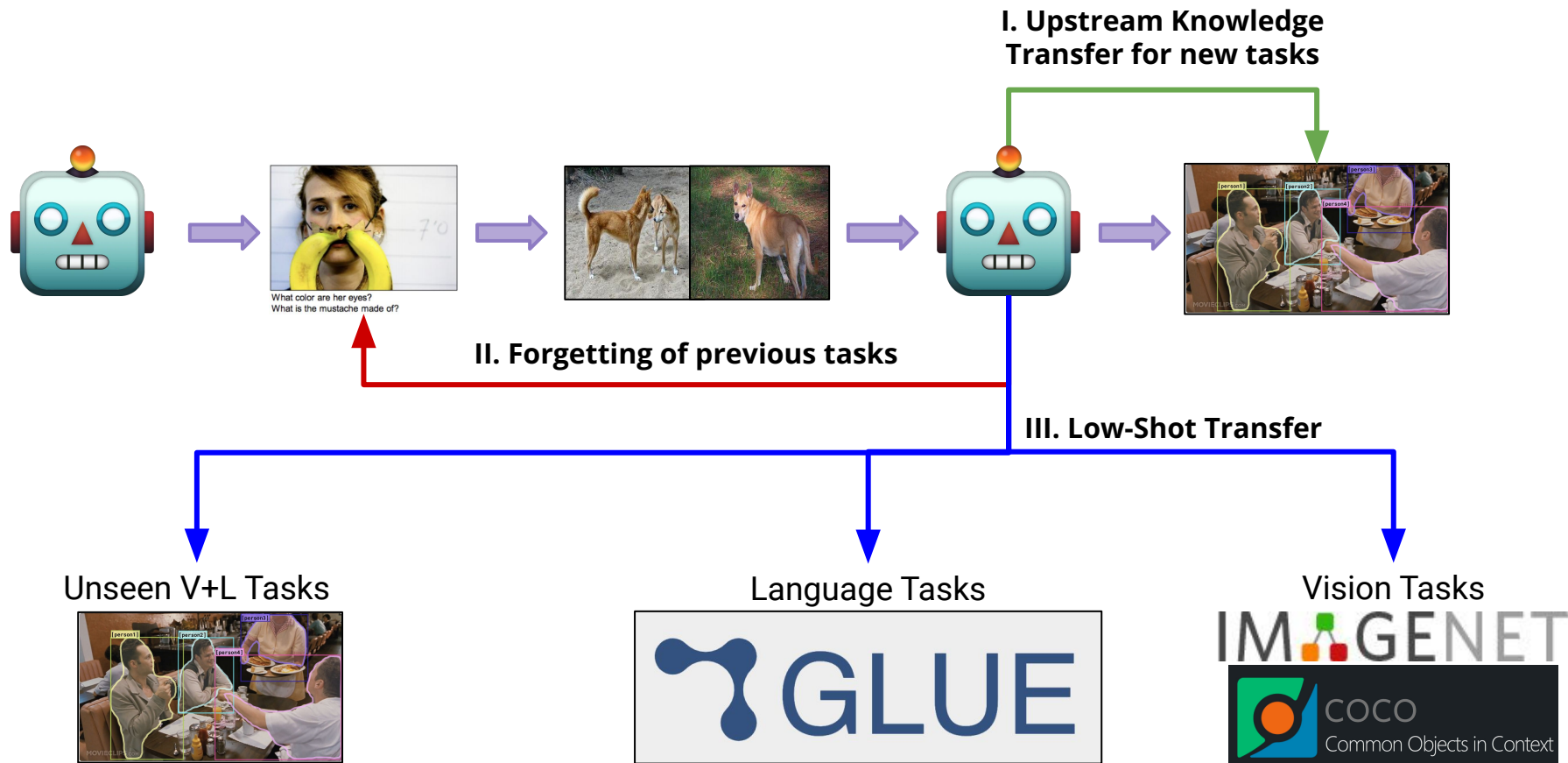
# Adapters

Insert new task-specific parameters into Transformer layers

- Transformer parameters kept frozen - **no forgetting!**
- **Fewer learnable parameters, faster to train**
- **Comparable performance as full model fine-tuning**
- **No cross-task knowledge transfer**



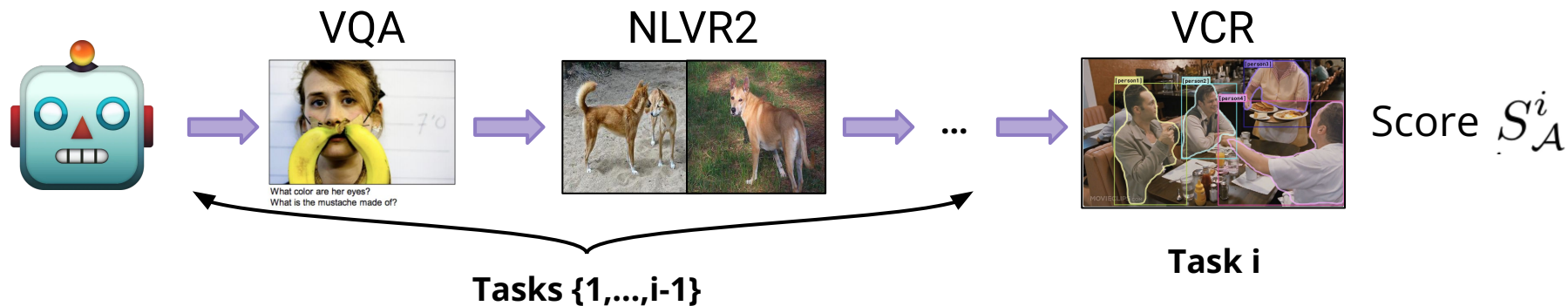
# IV. Evaluation



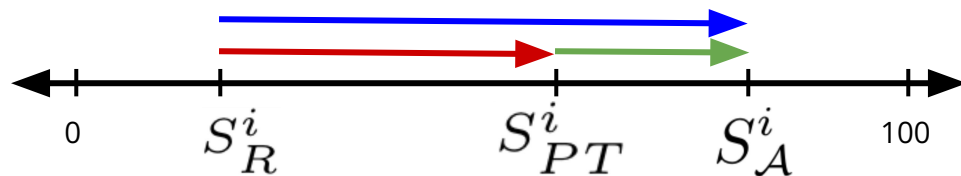
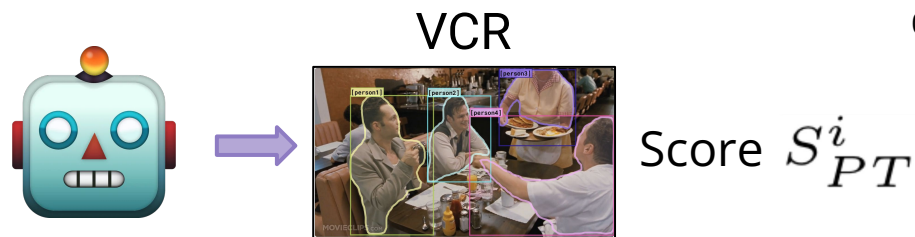


# Upstream Evaluation I: Upstream Knowledge Transfer

With Continual Learning Algorithm A:



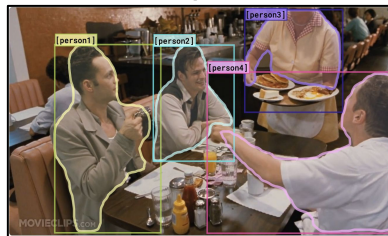
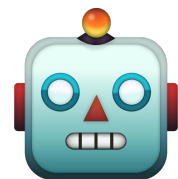
Without Continual Learning:



$$\mathbb{T}_{UK}(i) = \frac{S_A^i - S_{PT}^i}{S_{PT}^i - S_R^i}$$

# Upstream Evaluation II: Forgetting Transfer

Directly evaluate on previously learned task j



Score  $S_A^j$

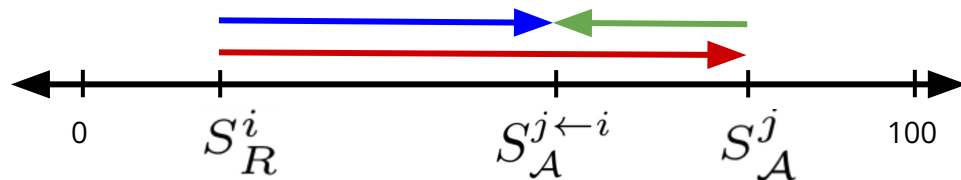


A man shoots a basketball  
**Contradiction**

Score  $S_A^{j \leftarrow i}$

Task j: VCR

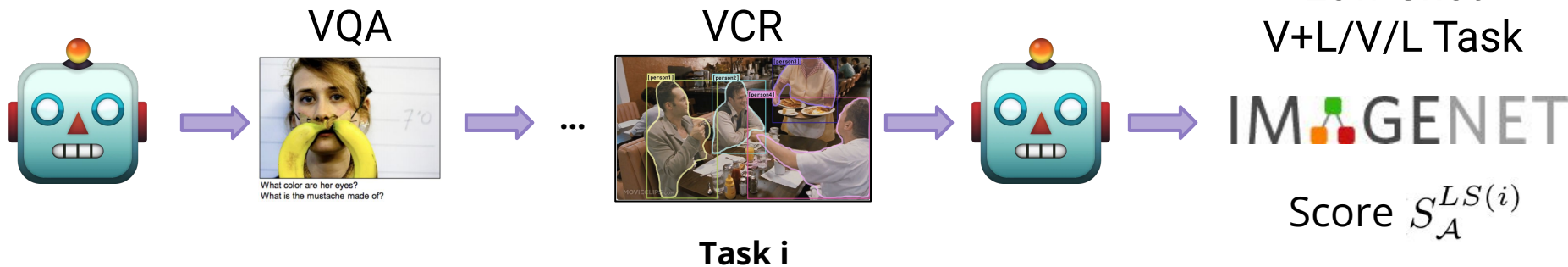
Task i: SNLI-VE



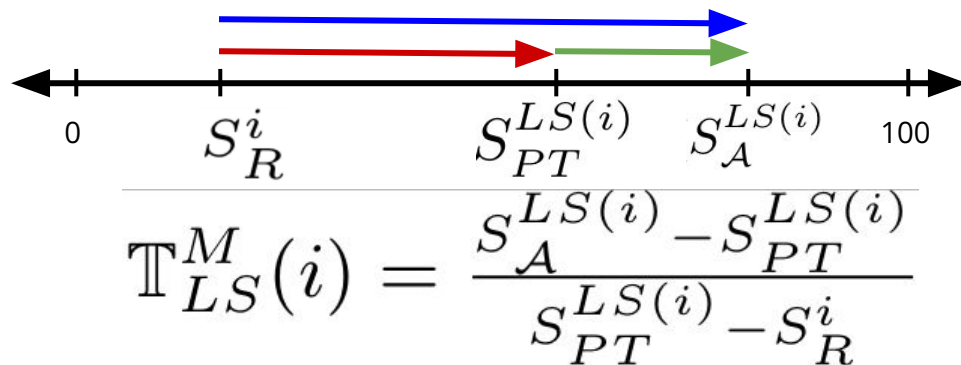
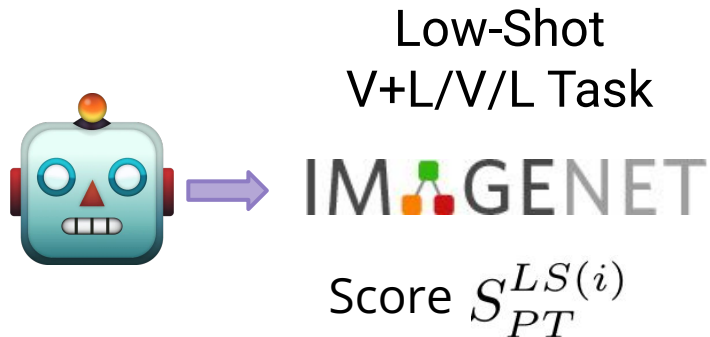
$$\mathbb{T}_F(j \leftarrow i) = \frac{S_A^j - S_A^{j \leftarrow i}}{S_A^j - S_R^i}$$

# Downstream Evaluation: Low-Shot Transfer

With Continual Learning Algorithm A:



Without Continual Learning:



# Experiments I: Upstream Continual Learning

- 4 V+L Tasks, ordered **VQA** → **NLVR2** → **SNLI-VE** → **VCR**
- **ViLT**-based continual learning model
- **6** different Continual Learning algorithms

# Results I: Upstream Continual Learning

**Upstream Knowledge Transfer:** How does Continual Learning affect model's ability to learn newly arriving tasks?

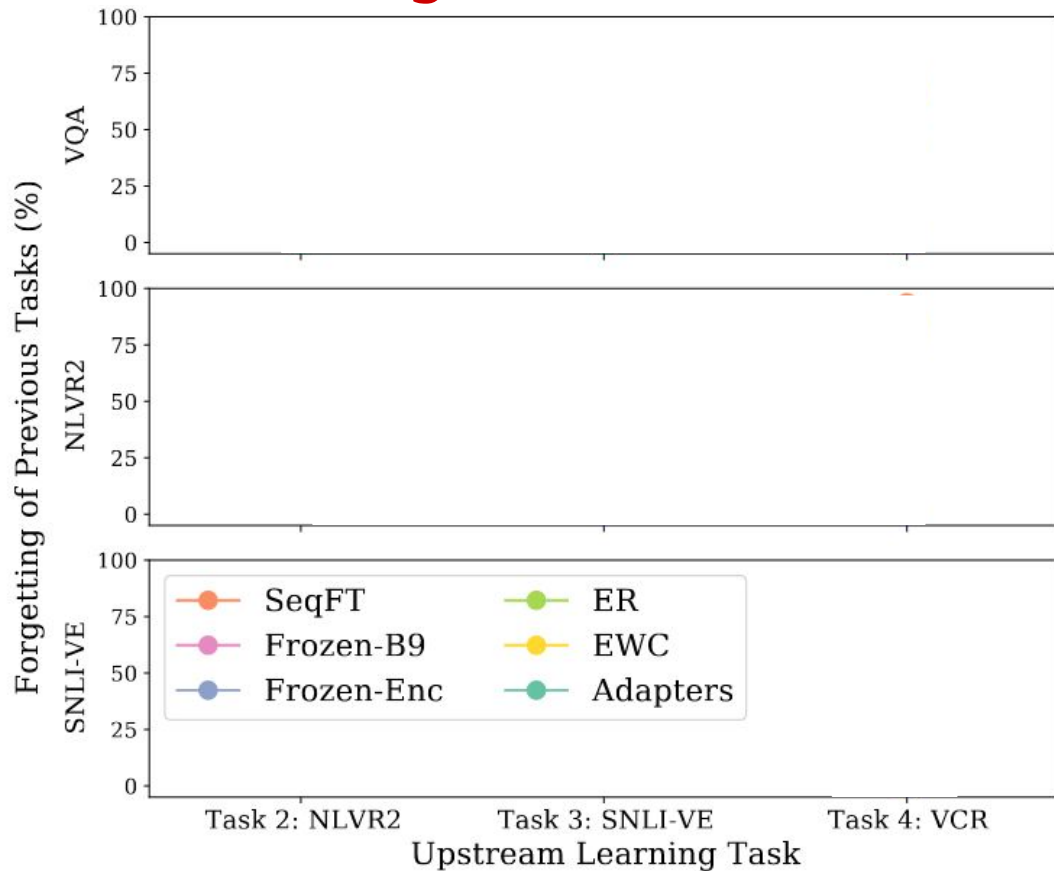
Alg $\mathcal{A}$	Params Trained	Task 1 VQAv2	Task 2 NLVR2	Task 3 SNLI-VE	Task 4 VCR
Direct FT	100%	[67.70]	[73.07]	[76.31]	[61.31]

- **More continual learning hurts ability to learn new tasks**
- **Adapters do not show negative transfer**, comparable to full model fine-tuning

# Results I: Upstream Continual Learning

**Forgetting**: How does learning new tasks affect model's performance on already-learned tasks?

- **More fine-tuned params == more forgetting**
- **ER > EWC**
- **Adapters >>>>**
- **Forgetting more severe after VCR**



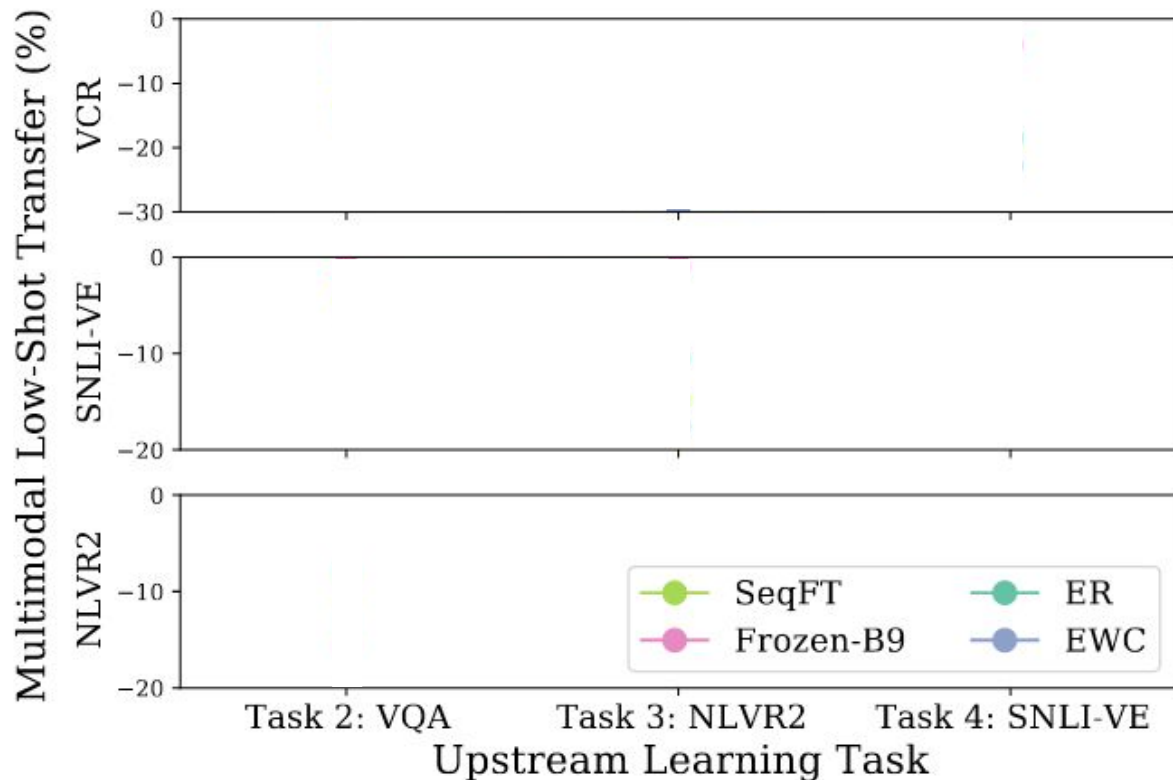
# Effect of Continual Learning Task Order

USC Viterbi School of Engineering

# Experiments and Results II: Downstream Low-Shot Transfer

Low-Shot Transfer to  
Unseen V+L Tasks

- Low-Shot transfer is always negative
- Unsurprising — CL also hurts model transfer with full training data



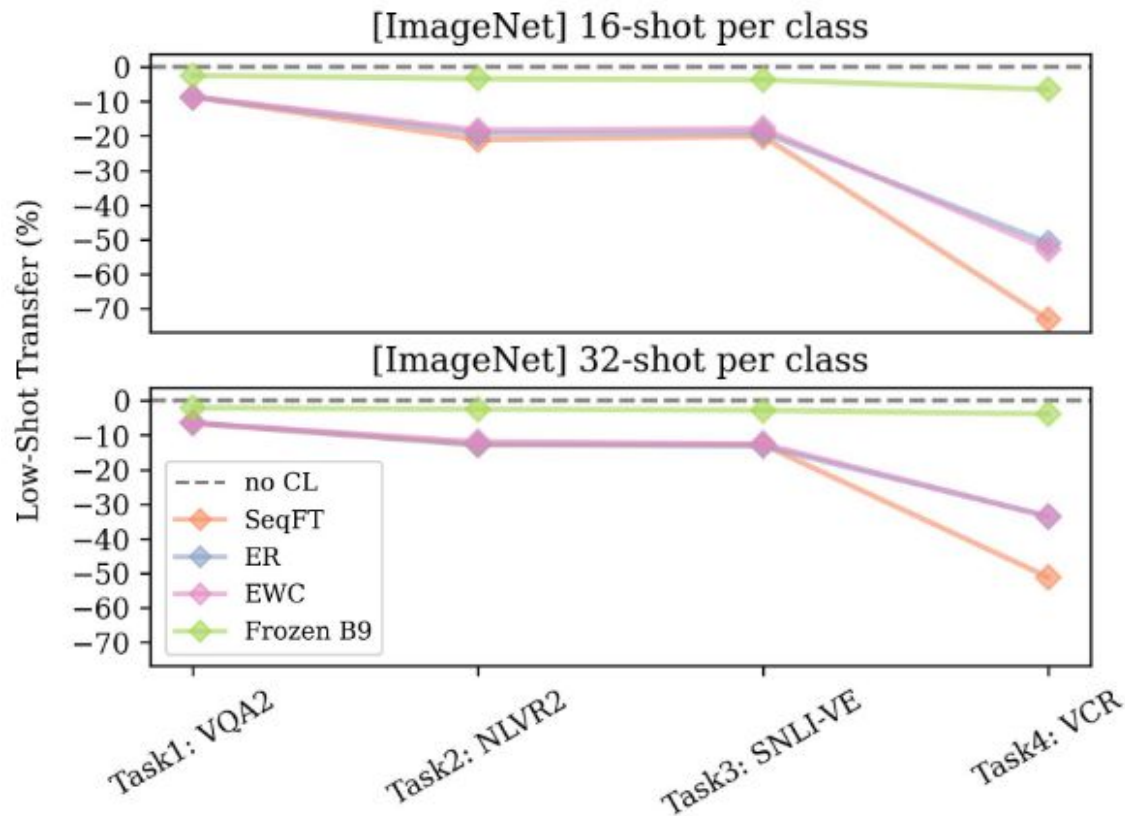


# Experiments and Results II: Downstream Low-Shot Transfer

## Low-Shot Transfer to Vision-Only Tasks

Language prompt: "This is an image."

- ViLT achieves good low-shot performance on vision tasks
- CL hurts low-shot transfer
- NLVR2 and VCR have more negative effect

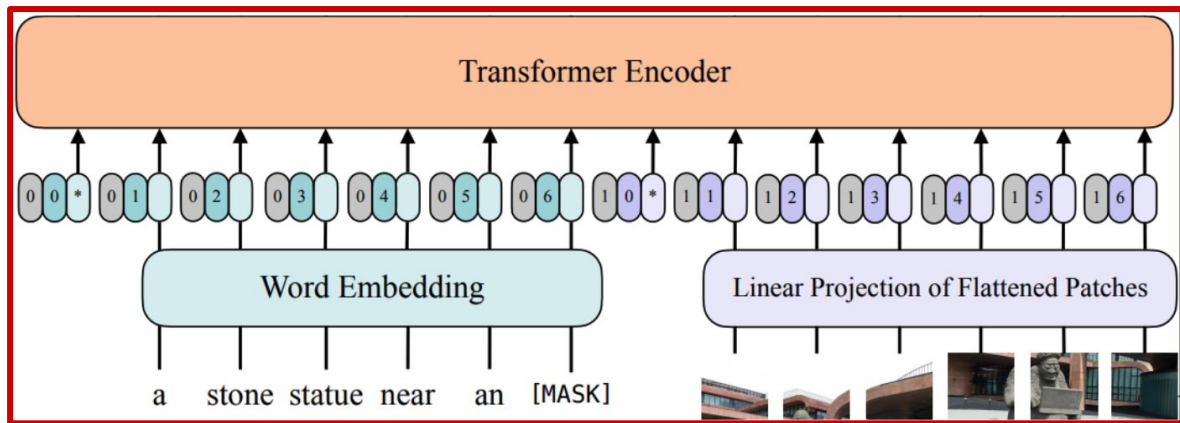


# Experiments and Results II: Downstream Low-Shot Transfer

## Low-Shot Transfer to Language-Only Tasks

### Adapting ViLT for NLP tasks:

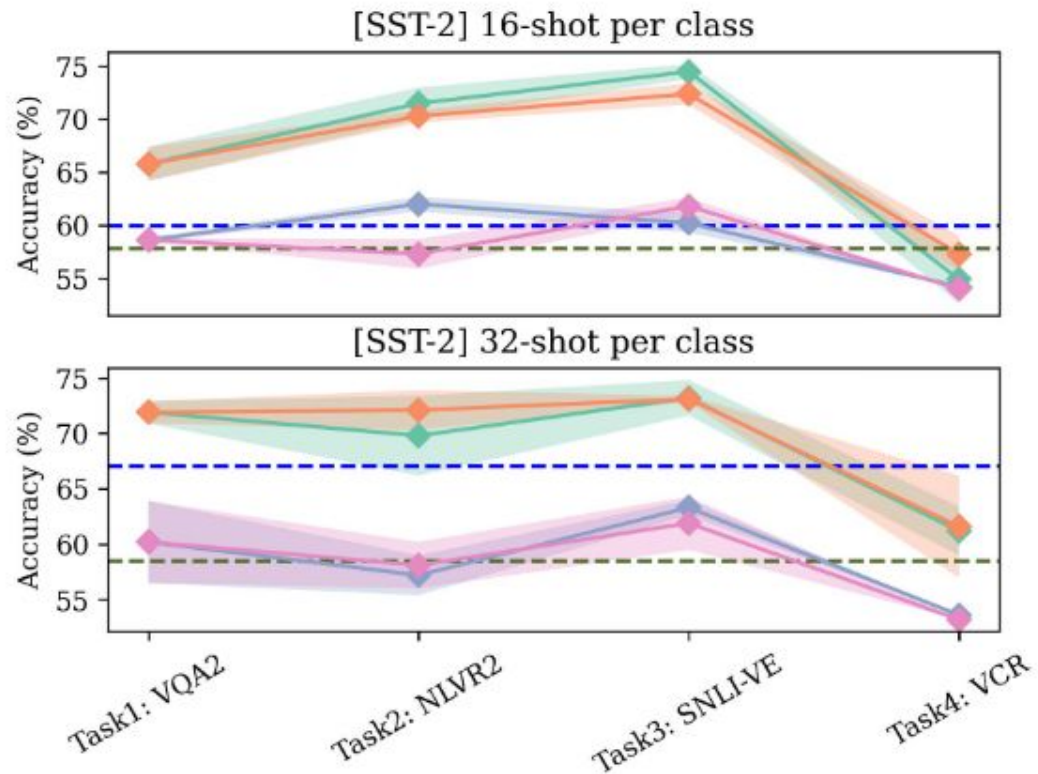
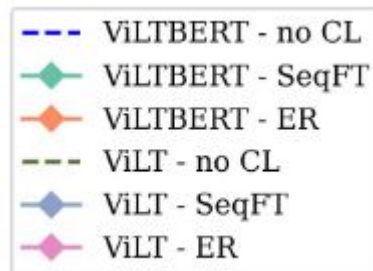
- Use **“average” MS-COCO image** for in-distribution visual input
- **Extend** language position embeddings
- **ViLT-BERT**: Replace language input embeddings with BERT representations



# Experiments and Results II: Downstream Low-Shot Transfer

## Low-Shot Transfer to Language-Only Tasks

- Upstream CL helps! Sometimes
- ViLT sees negligible differences
- CL helps ViLT-BERT with SST2
- VCR hurts SST2
- CL hurts multi-choice tasks



# Conclusions

- We propose **CLiMB**, a benchmark to study CL in multimodal settings
- CLiMB is an **extensible community tool** for studying tasks, model architectures, and CL algorithms.
- **Existing Continual Learning methods fail** at:
  - generalizing well to sequences of multimodal tasks
  - Enabling low-shot adaptation to multi/unimodal tasks
- **Adapters are most effective** at preserving pre-trained model knowledge and forgetting mitigating
- There is **a need for new research** into Continual Learning strategies for this challenging multimodal setting.

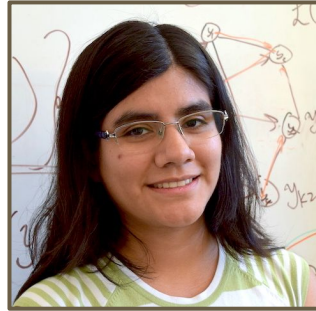
# Future Directions

- **Adapters that share knowledge across tasks**
- **Multimodal Adapters**
- Studying **multimodal distribution shifts**
- Building a **task-agnostic** modeling framework:
  - Sequence-to-sequence task formulations
  - Integrating **generalist models** into CLiMB
  - **Embodied** navigation, task completion

# Acknowledgements



Ting-Yun Chang ✨



Leticia Pinto Alva ✨



Georgios Chochlakis



Mohammad Rostami



Jesse Thomason ✨

# Thank You!!

<https://github.com/GLAMOR-USC/CLiMB>

GLAMOR-USC / CLiMB (Public)

<> Code Issues Pull requests Actions Projects Wiki Security Insights Settings

main 1 branch 0 tags Go to file Add file Code

tejas1995 Delete plot\_forgetting.py beee722 15 days ago 213 commits

figs	Code update 8/15	15 days ago
src	Delete plot_forgetting.py	15 days ago
.gitignore	Added wandb to gitignore	6 months ago
.gitmodules	Added ViitModel to adapter-transformers repo fork	5 months ago
ADD_NEW_ALGORITHMS.md	Code update 8/15	15 days ago
ADD_NEW_MODELS.md	Code update 8/15	15 days ago
ADD_NEW_TASKS.md	Code update 8/15	15 days ago
DATA_DOWNLOAD.md	Code update 6/20/22	2 months ago
LICENSE	Update LICENSE	2 months ago
README.md	Code update 8/15	15 days ago
TRAIN_UPSTREAM_CL.md	Code update 8/15	15 days ago
requirements.txt	Training scripts, commented out transformers in requirements	3 months ago

README.md

## CLiMB: The Continual Learning in Multimodality Benchmark