

Multimodal AI

Lecture 10.2 – Cross-modal Transfer

Paul Liang

Assistant Professor

MIT Media Lab & MIT EECS



<https://pliang279.github.io>

ppliang@mit.edu

 [@pliang279](#)



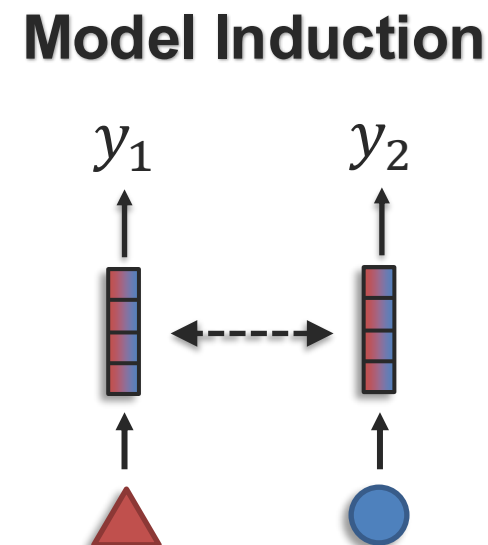
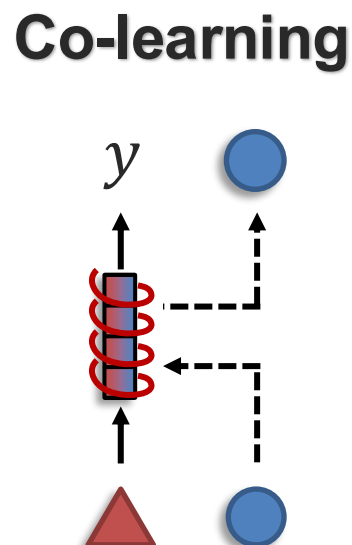
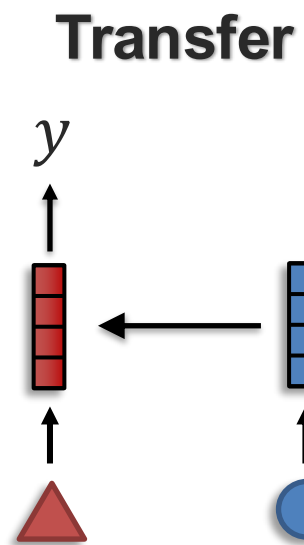
Today's lecture

- 1 Basics of cross-modal transfer
- 2 Cross-modal transfer via fusion
- 3 Cross-modal transfer via alignment
- 4 Cross-modal transfer via translation

Transference

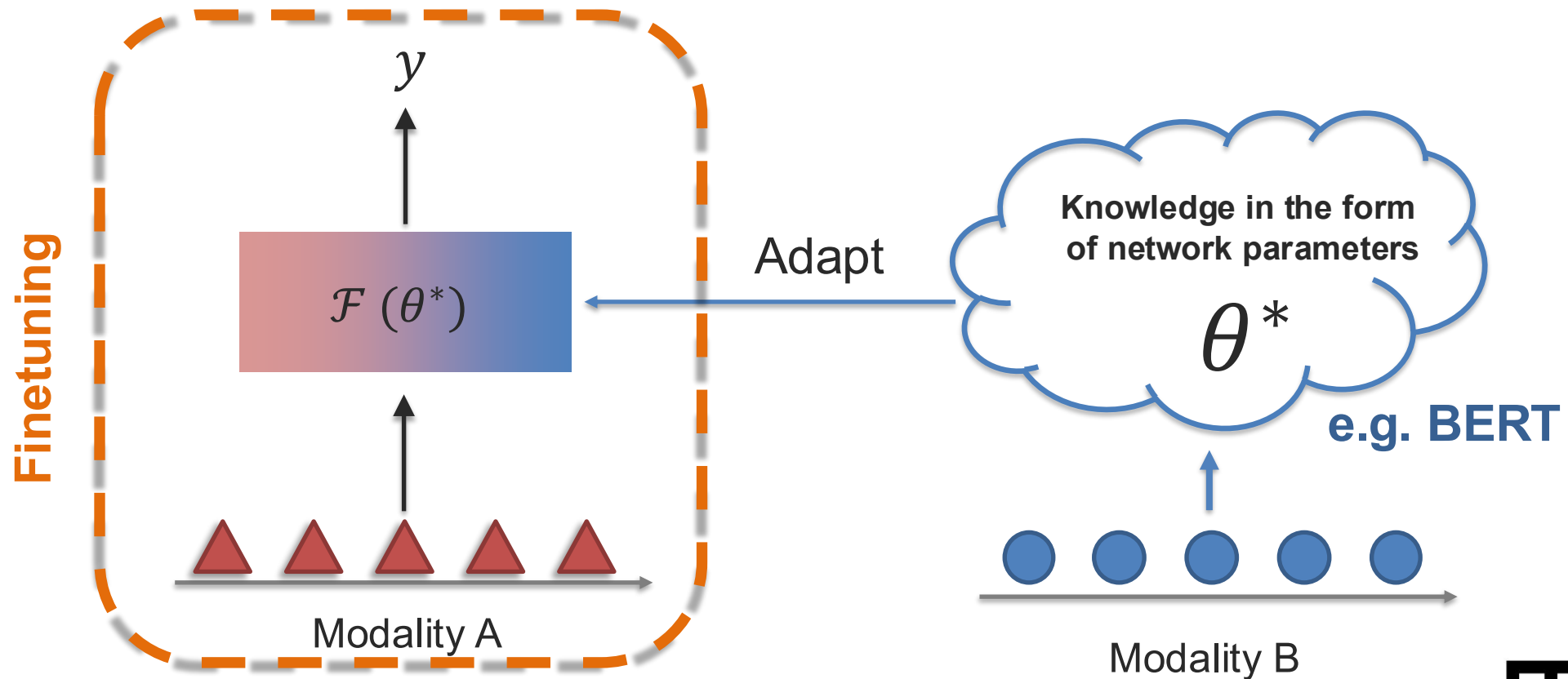
Definition: Transfer knowledge between modalities, usually to help the primary modality which may be noisy or with limited resources

Sub-challenges:



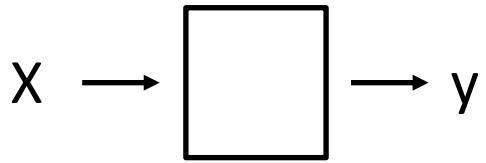
Part 1: Transfer via Pretrained Models

Definition: Transferring knowledge from large-scale pretrained models to downstream tasks involving the primary modality.

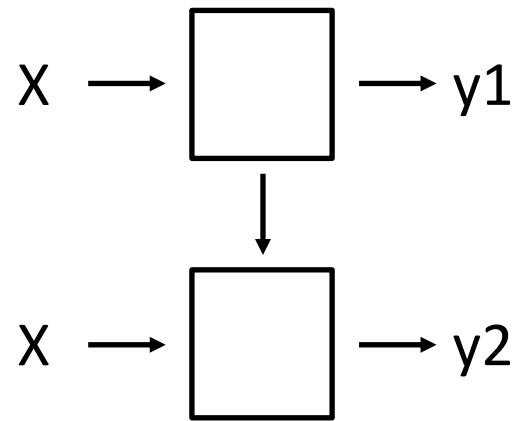


Multitask and Transfer Learning

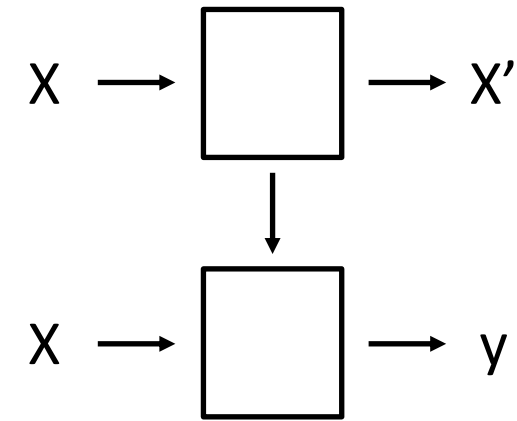
Supervised learning



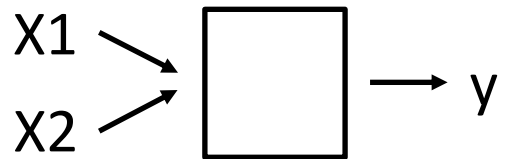
Transfer learning



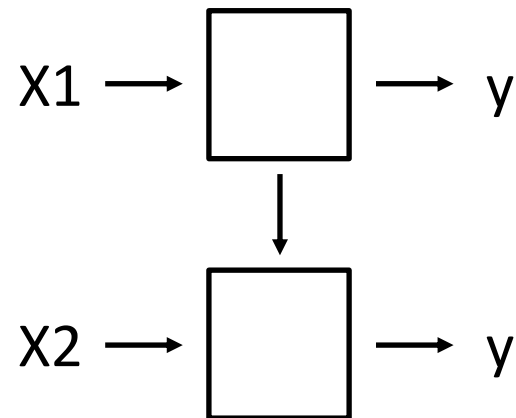
Unsupervised/self-supervised pre-training



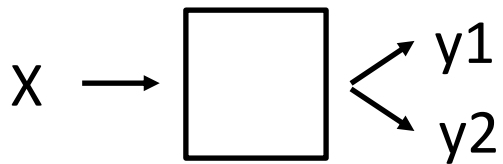
Multimodal (supervised) learning



Cross-modal learning


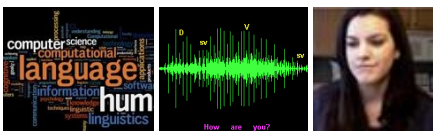


Multitask learning



Multitask and Transfer Learning

Humans

Language Speech Gestures

Multimedia




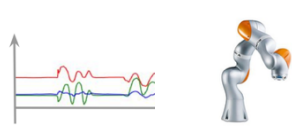



Image Language

Robotics





Force Proprioception

Healthcare

Design

Generalization across modalities and tasks
Important if some tasks are low-resource



SUBJECT_ID

Age
Sex
Ethnicity
...

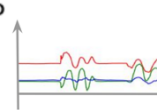


Table Sensors



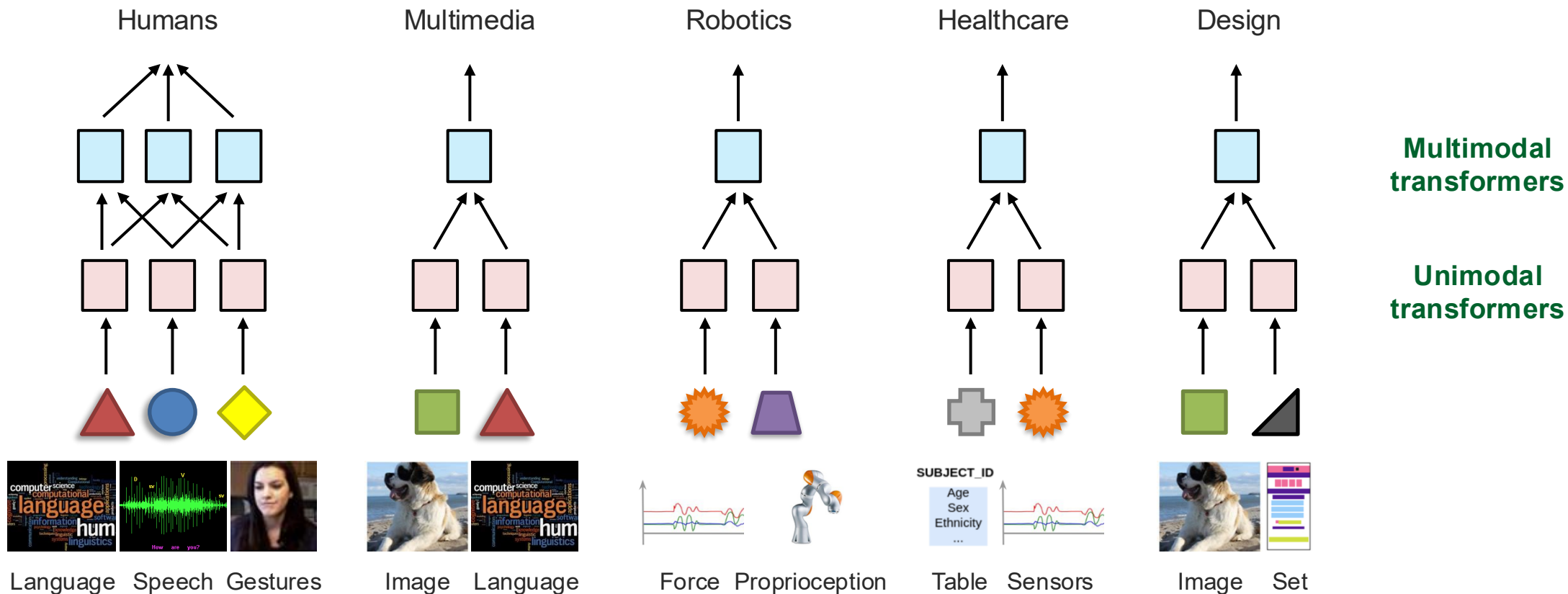



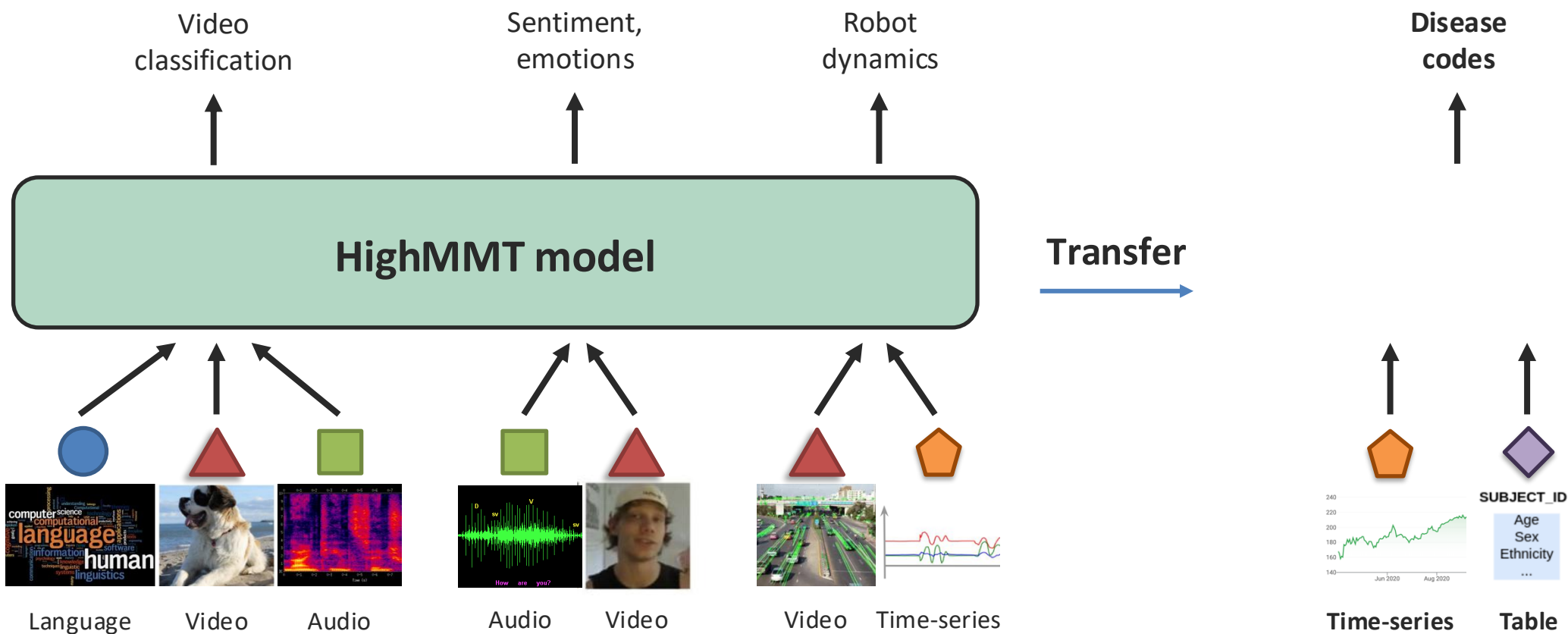
Image Set

Multitask and Transfer Learning



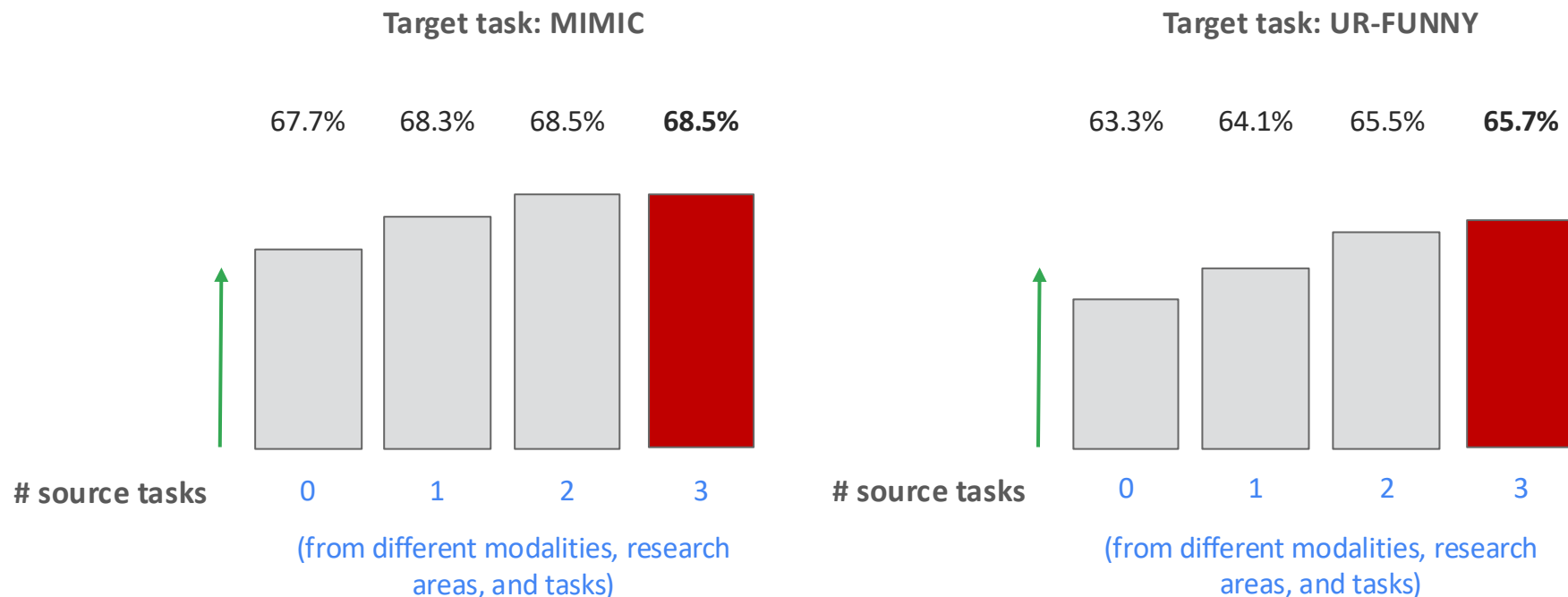
High-Modality Multimodal Transformer

Transfer across partially observable modalities



Multitask and Transfer Learning

Transfer across partially observable modalities

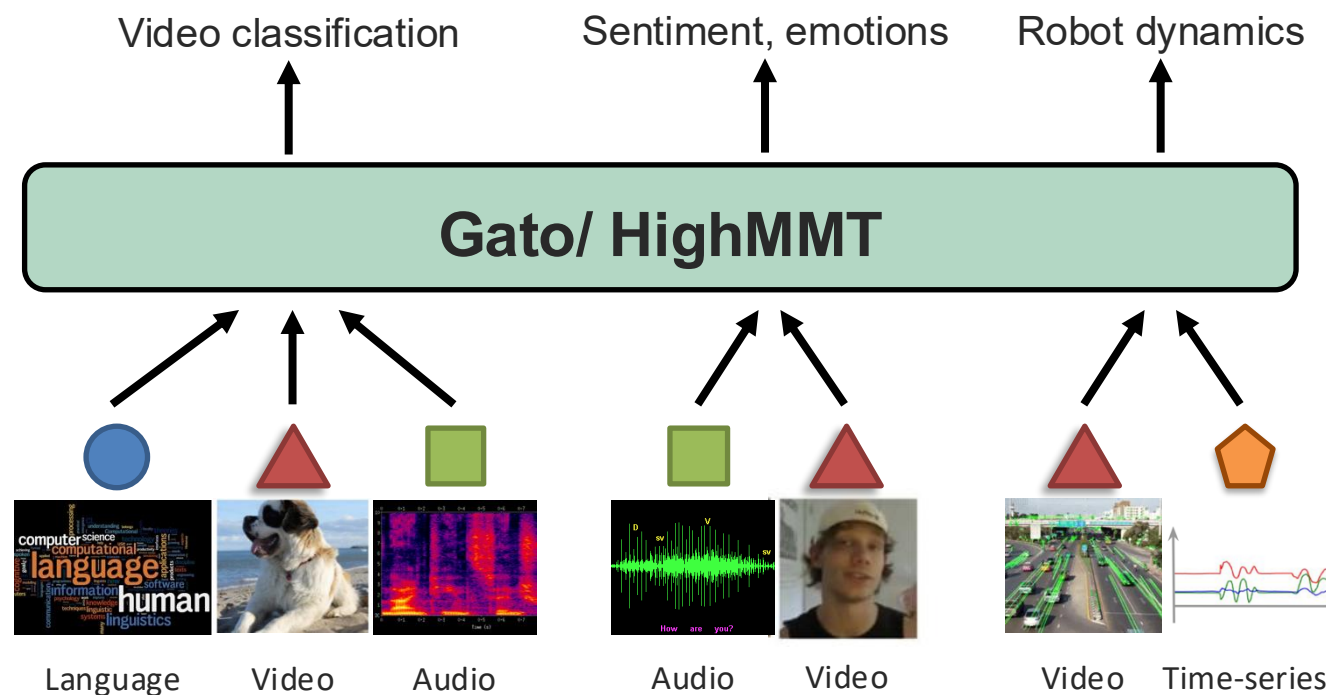


Achieves both multitask and transfer capabilities across modalities and tasks

High-Modality Models

Some implicit assumptions:

- All modalities can be represented as sequences without losing information.
- Dimensions of heterogeneity can be perfectly captured by modality-specific embeddings.
- Cross-modal connections & interactions are shared across modalities and tasks.



Shared multimodal model?

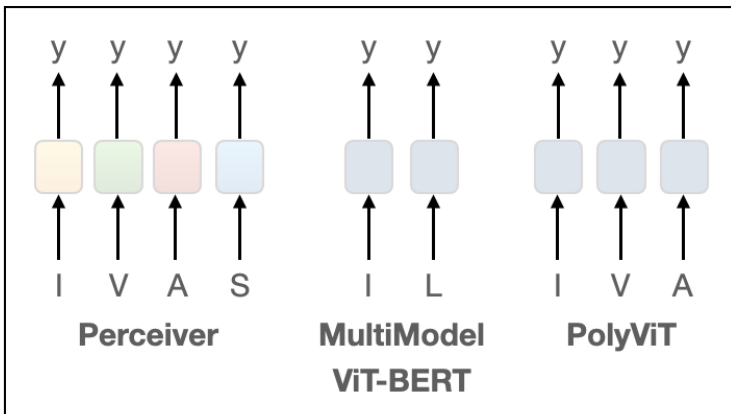
Modality-specific embeddings?

Standardized input sequence?

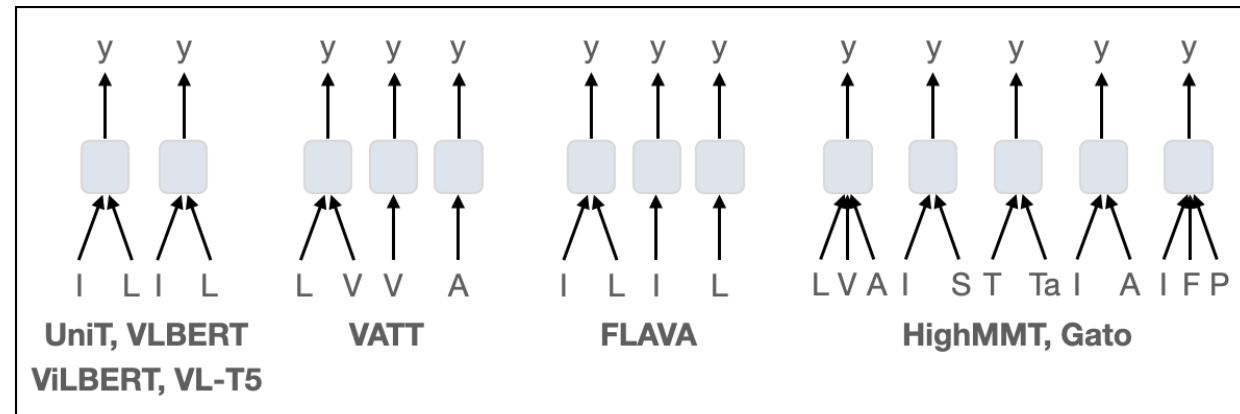
Multitask and Transfer Learning

Many more dimensions of transfer

Unified encoder for unimodal learning



Multimodal multitask learning



I: image
 V: video
 A: audio
 S: set
 L: language
 T: time-series
 Ta: tables
 F: force sensor
 P: proprioception sensor

common architecture

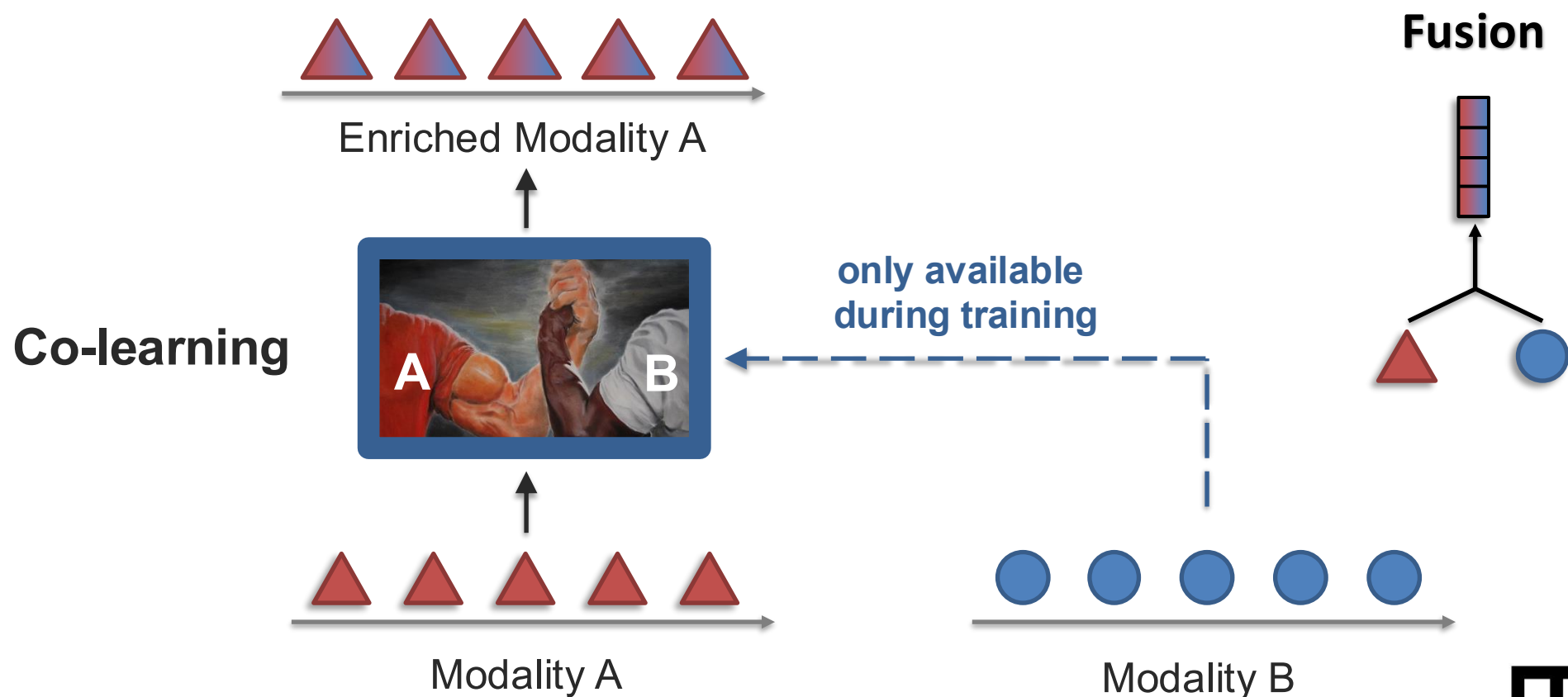
parameter sharing

Open challenges:

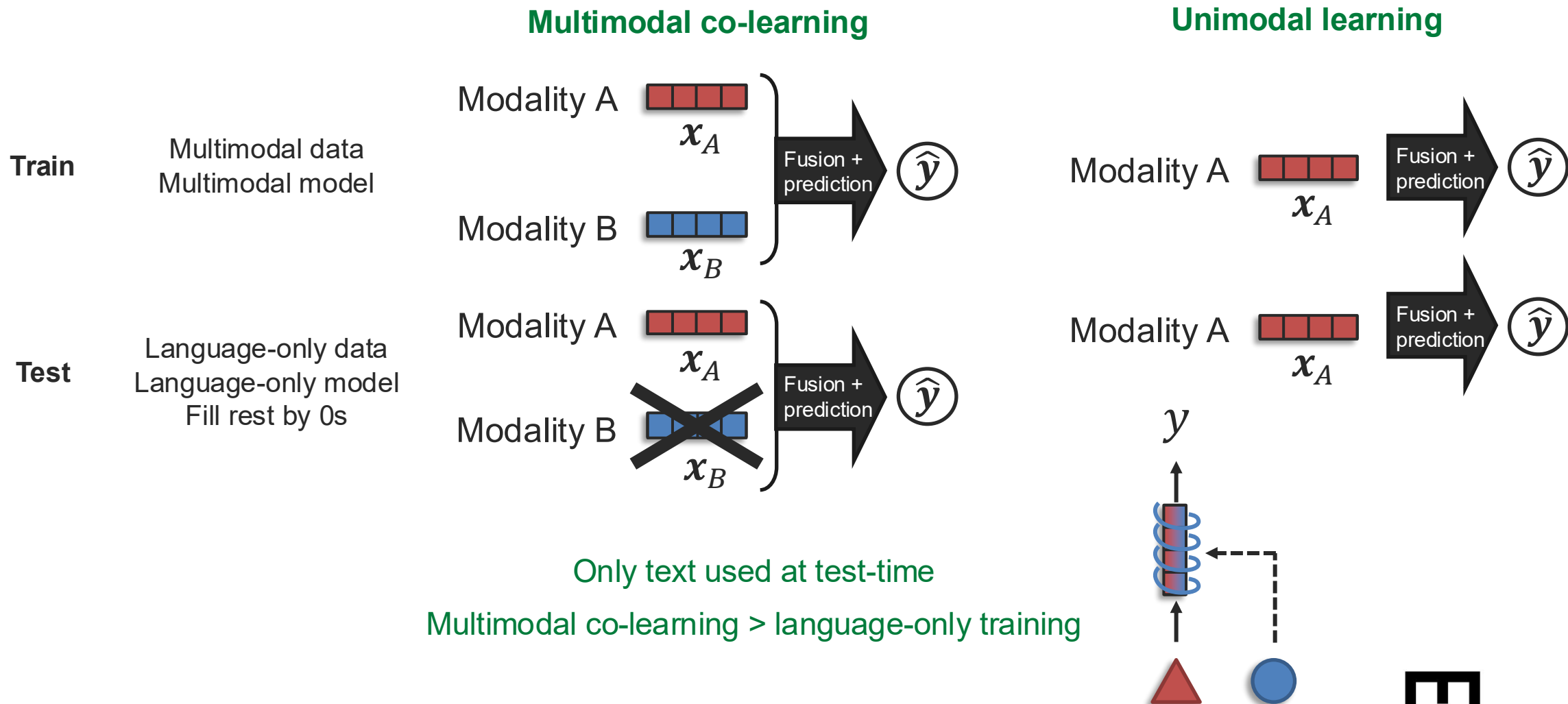
- Low-resource: little downstream data, lack of paired data, robustness (next section)
- Beyond language and vision
- Settings where SOTA unimodal encoders are not deep learning e.g., tabular data
- Complexity in data, modeling, and training
- Interpretability (next section)

Part 2: Co-learning

Definition: Transferring information from secondary to primary modality by sharing representation spaces between both modalities.



Co-learning via Fusion



Co-learning via Fusion

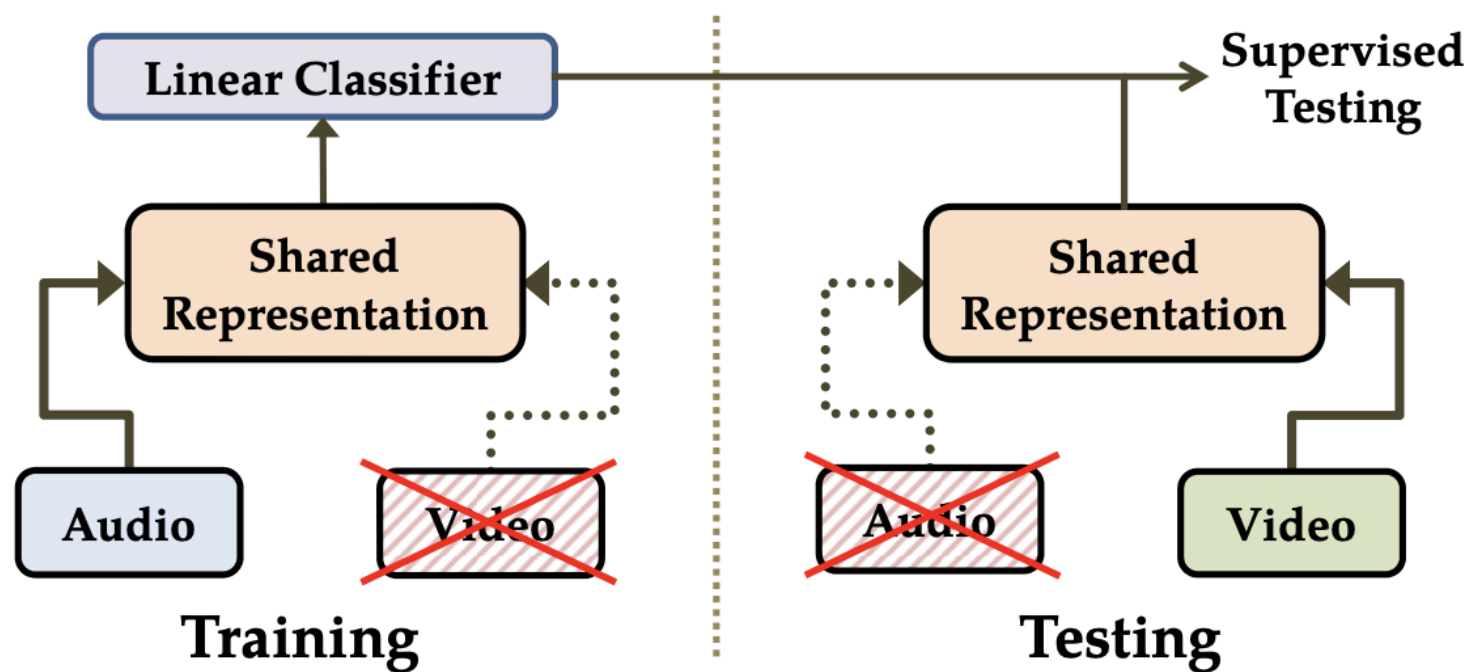
Generative model (Deep Boltzmann Machine) to learn joint representation and infer missing text.

Model	MAP	Prec@50
Image LDA (Huiskes et al., 2010)	0.315	-
Image SVM (Huiskes et al., 2010)	0.375	-
Image DBN	0.463 \pm 0.004	0.801 \pm 0.005
Image DBM	0.469 \pm 0.005	0.803 \pm 0.005
Multimodal DBM (generated text)	0.531 \pm 0.005	0.832 \pm 0.004

learning multimodal features helps even when some modalities are absent at test time.

Co-learning via Fusion

Train on some subset of modalities and test on another

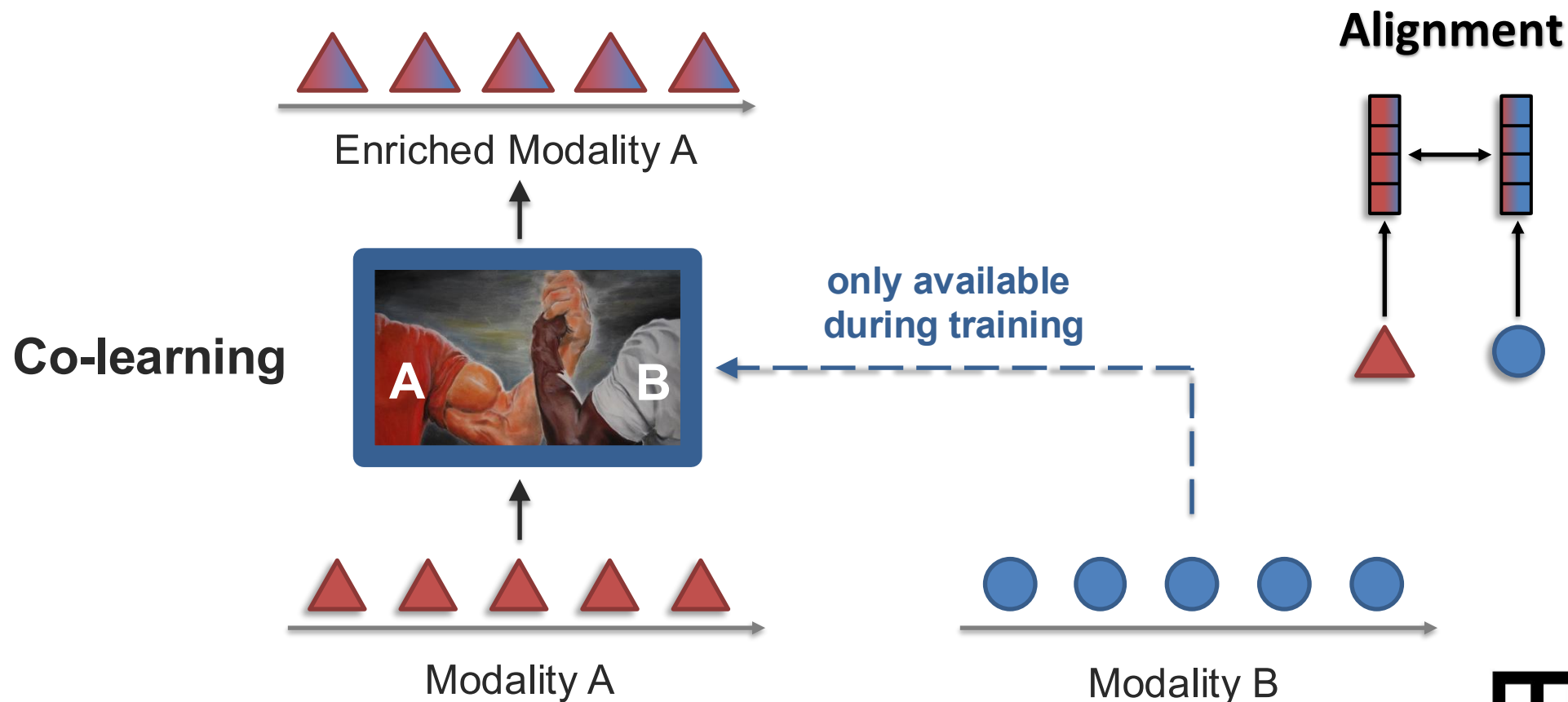


McGurk effect



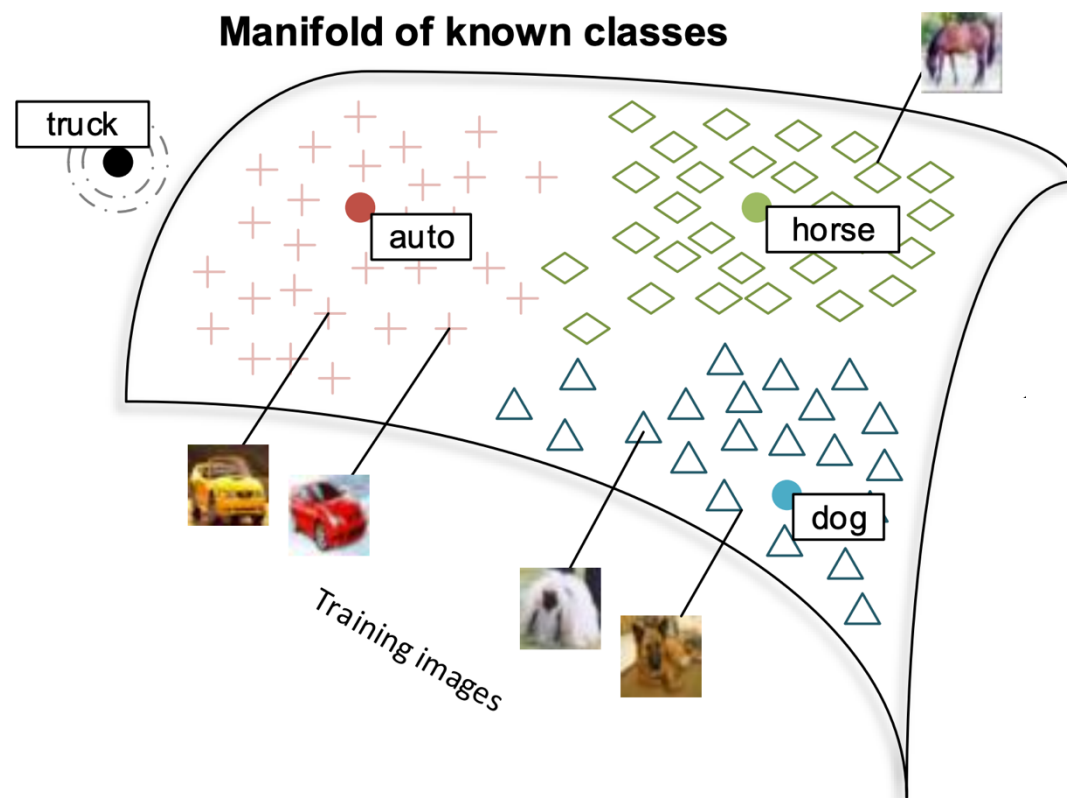
Co-learning via Alignment

Definition: Transferring information from secondary to primary modality by sharing representation spaces between both modalities.

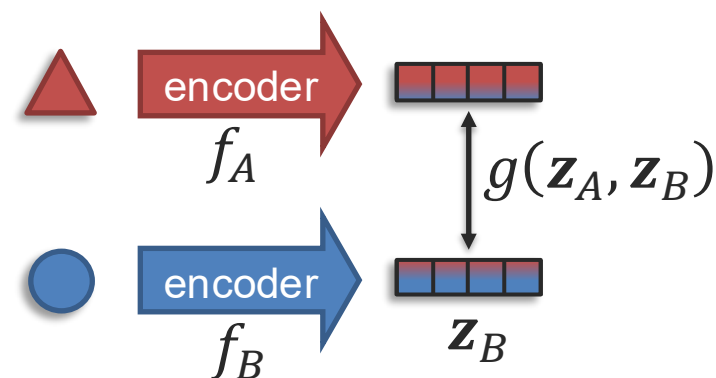


Co-learning via Alignment

Representation alignment: word embedding space for zero-shot visual classification

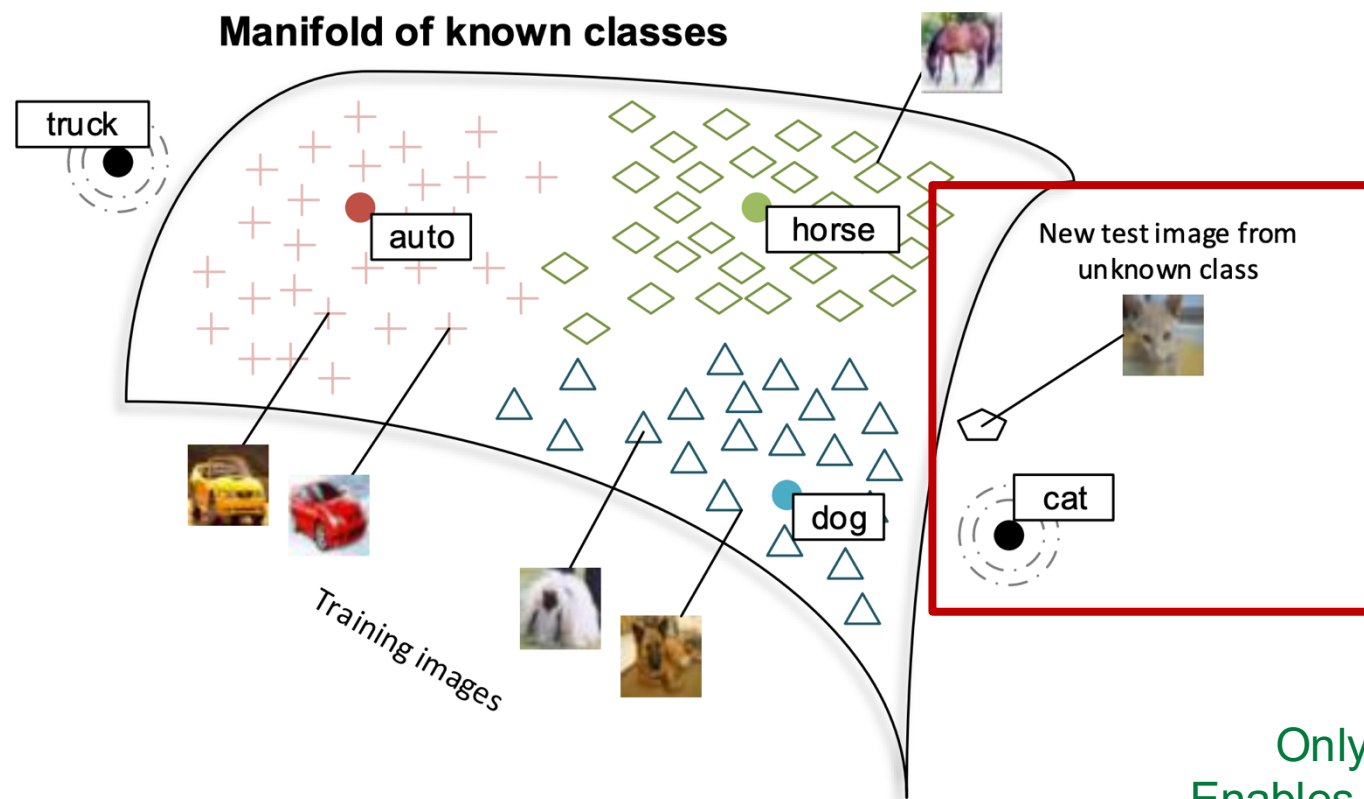


Recall representation alignment!

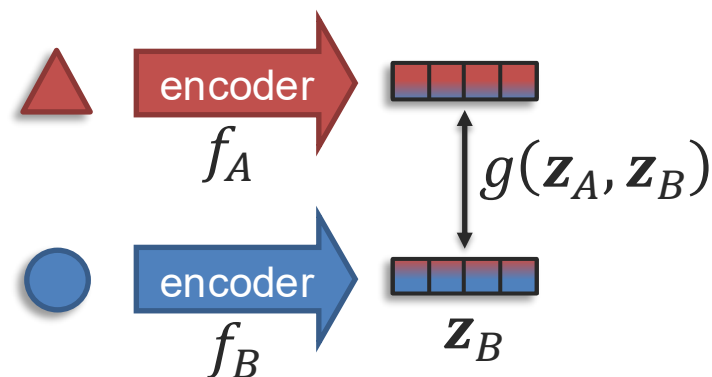


Co-learning via Alignment

Representation alignment: word embedding space for zero-shot visual classification



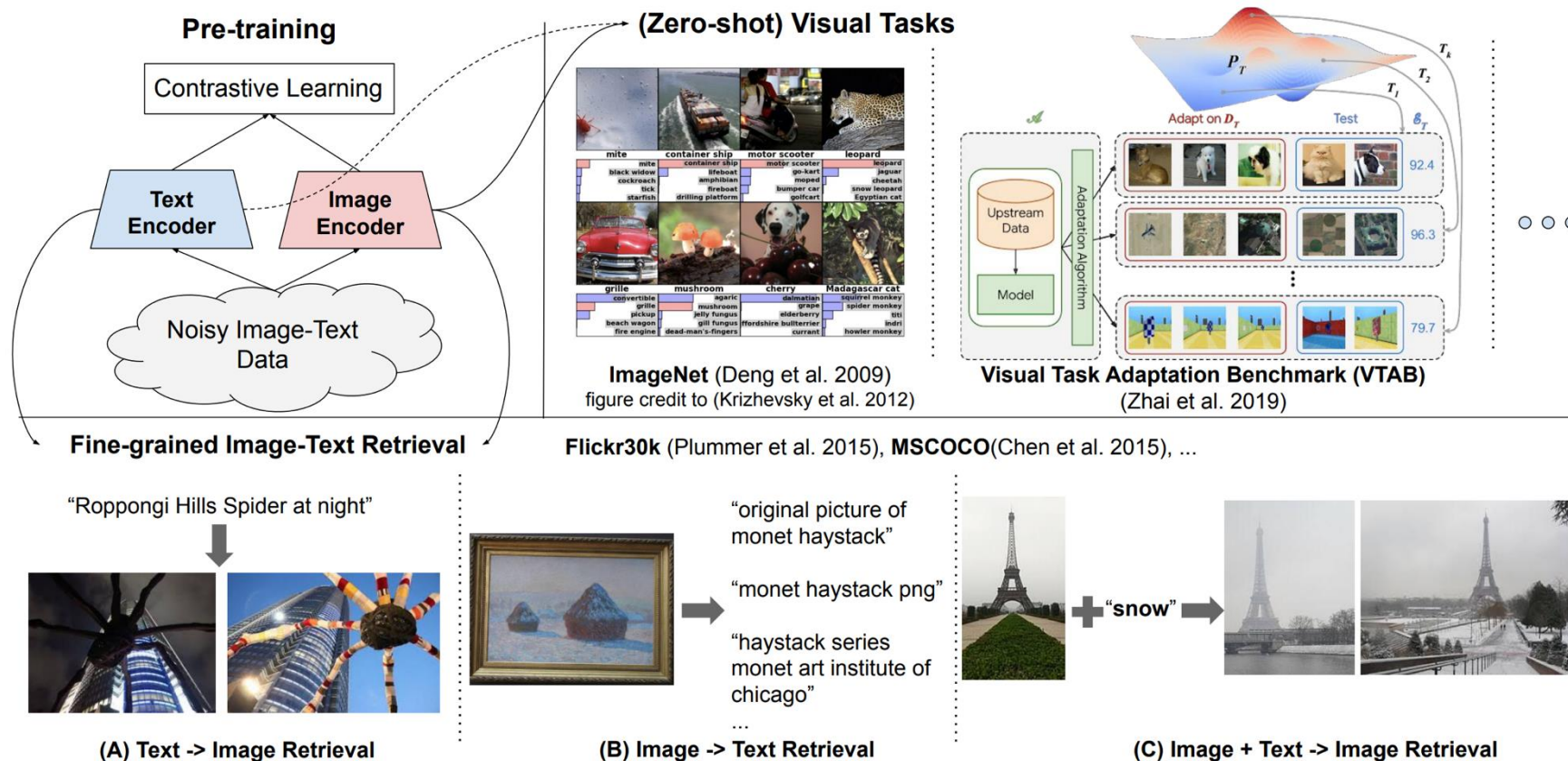
Recall representation alignment!



Only images used at test-time
Enables zero-shot image classification

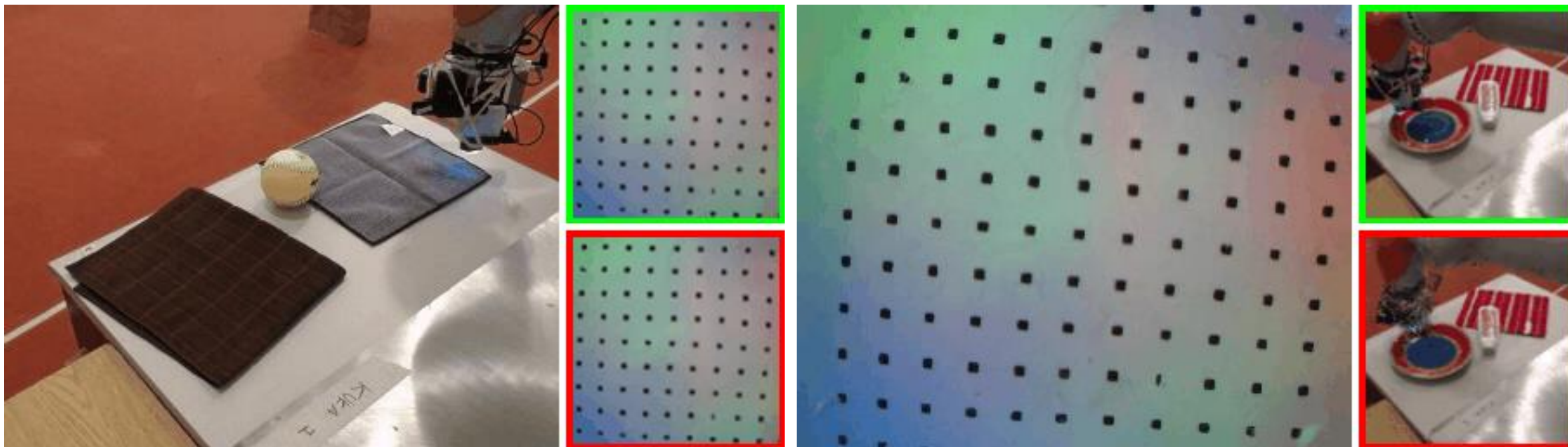
Co-learning via Alignment

Representation alignment at scale



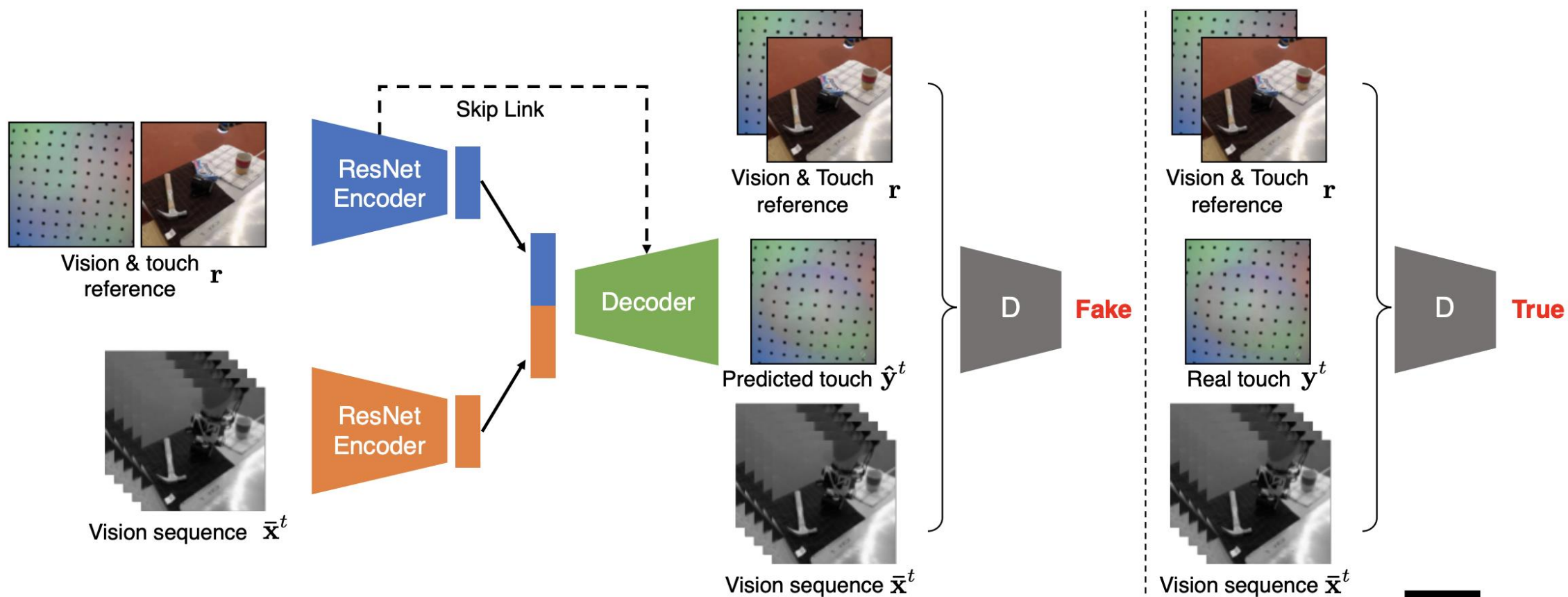
Vision-Touch Alignment

Aligning vision and touch in robotics



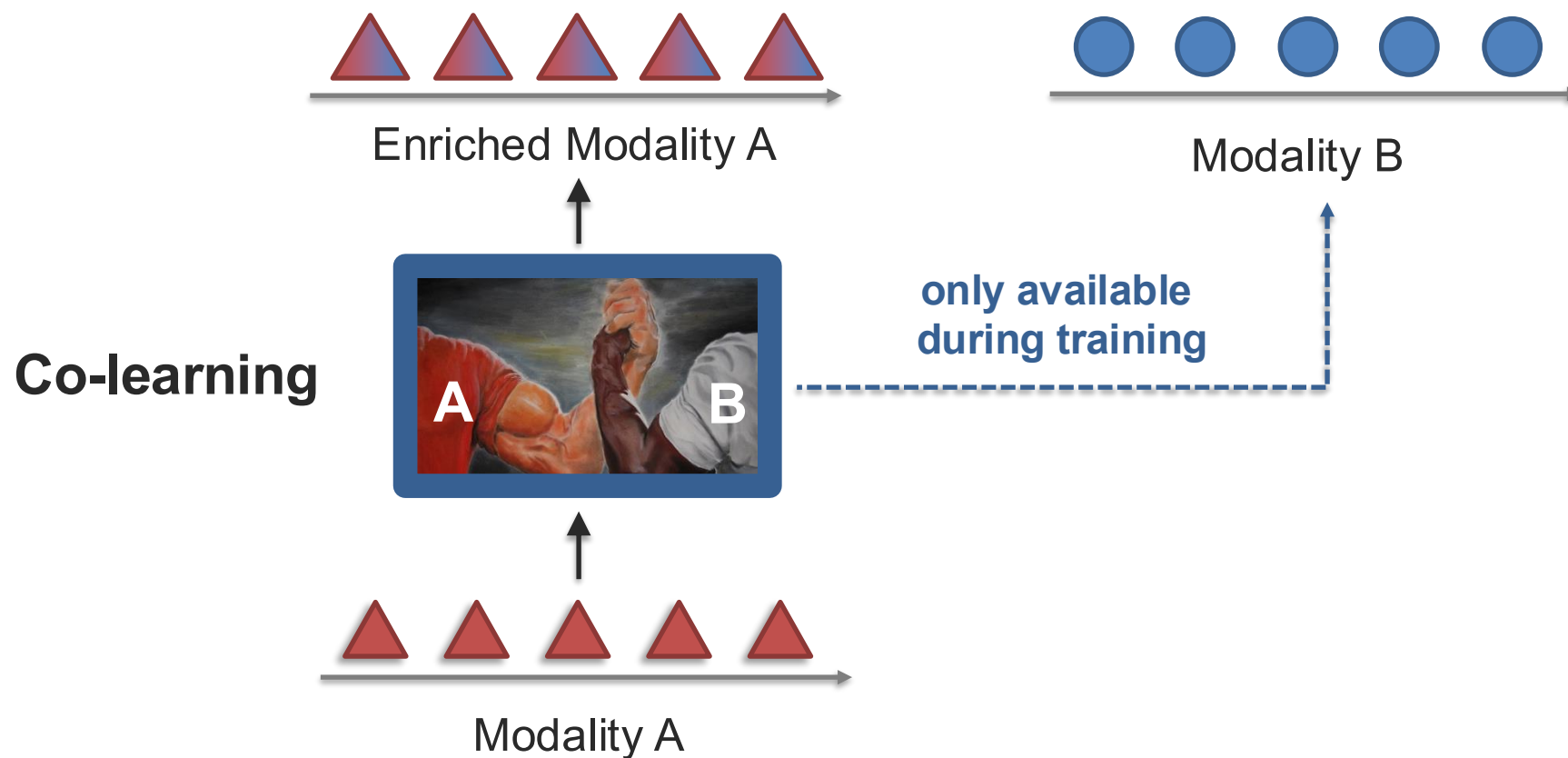
Vision-Touch Alignment

Aligning vision and touch in robotics



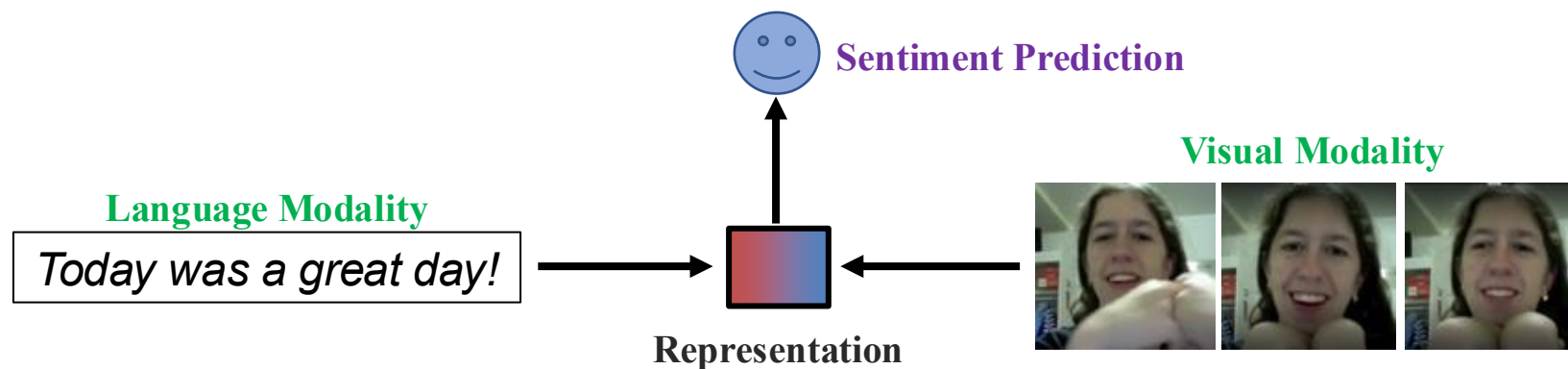
Co-learning via Translation

Definition: Transferring information from secondary to primary modality by using the secondary modality as a generation target.



Co-learning via Translation

Bimodal translations

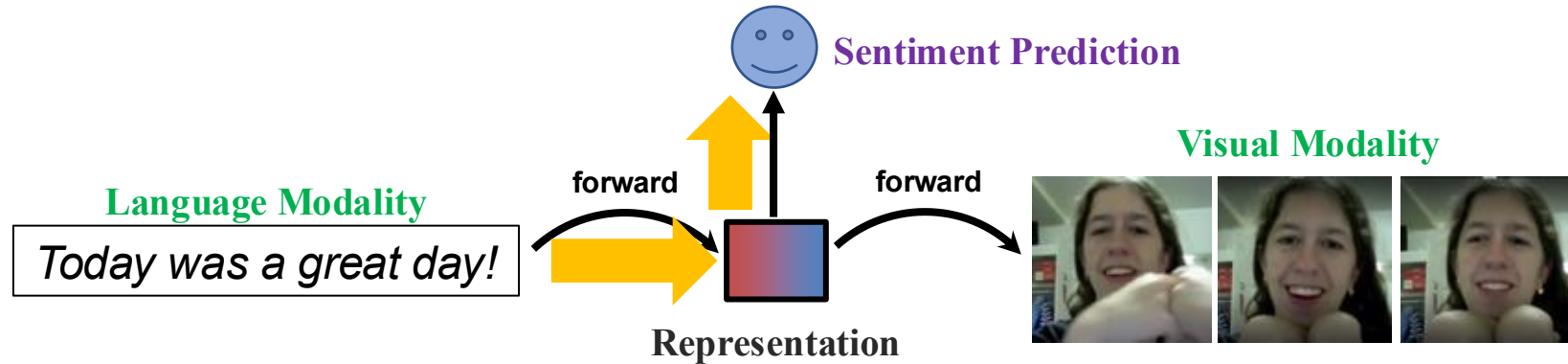


Both modalities required at test time!
Sensitive to noisy/missing visual modality.

We want to leverage information from visual modality
while being robust to it during test-time.

Co-learning via Translation

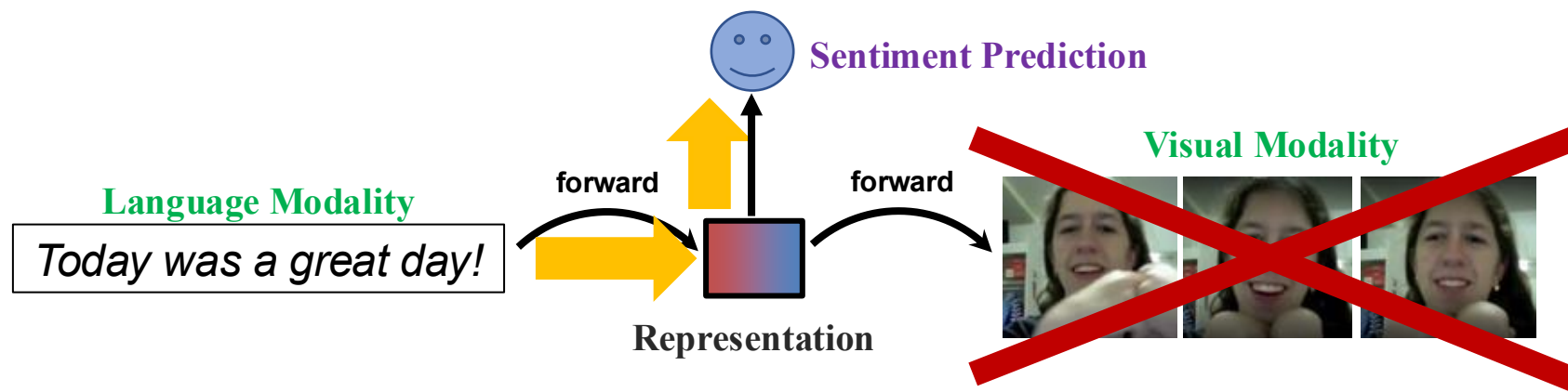
Bimodal translations



Cross-modal translation during training
Only language modality required at test time!

Co-learning via Translation

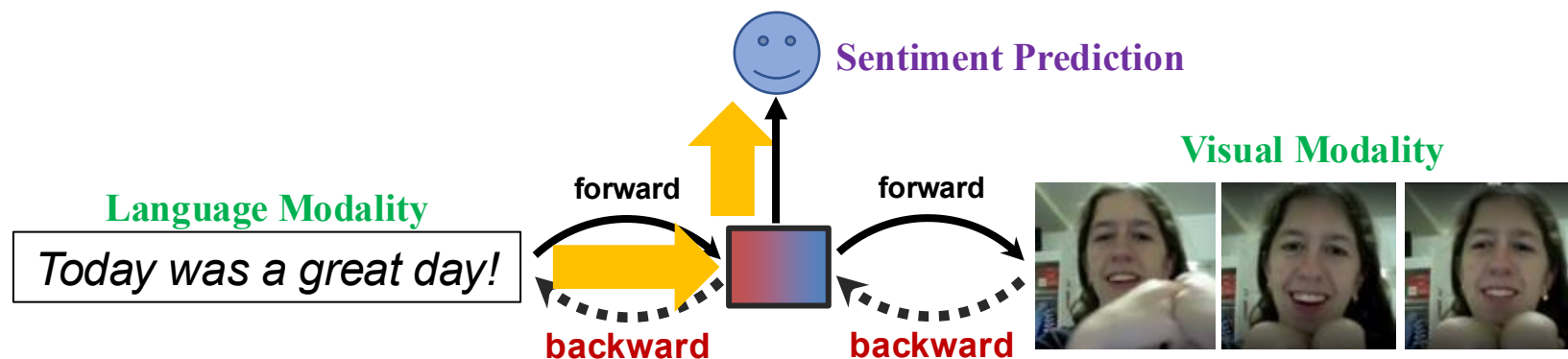
Bimodal translations



Problem: how do you ensure that both modalities are being used?

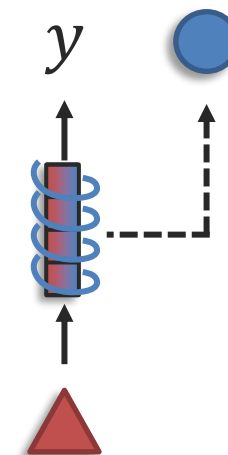
Co-learning via Translation

Bimodal cyclic translations



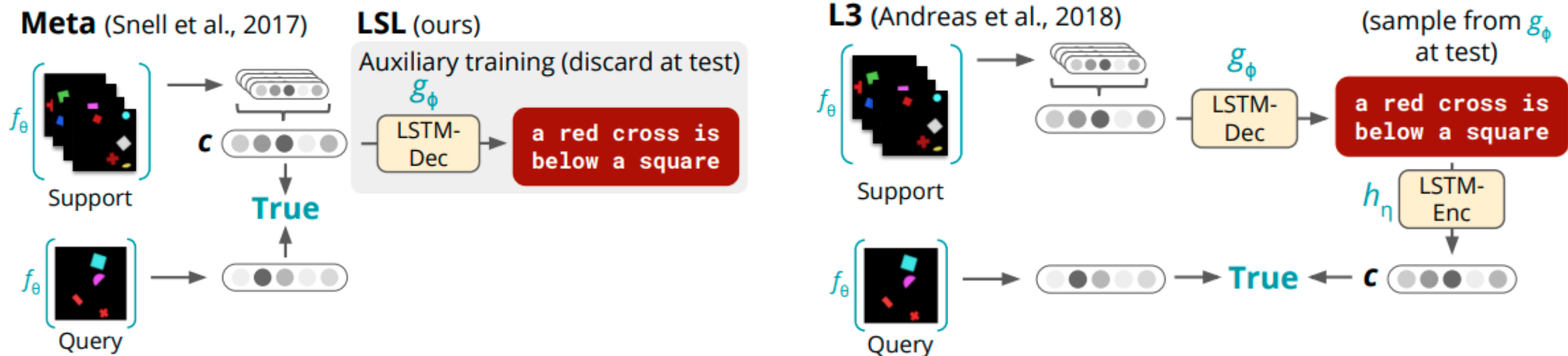
Solution: cyclic translations from visual back to language

Cross-modal translation during training
Only language modality required at test time!



Co-learning for Compositionality

Image to text translation



[Mu et al., 2019. Shaping Visual Representations with Language for Few-Shot Classification]

[Andreas et al. 2017, Learning with Latent Language]

[Sharma et al. 2021. Skill Induction and Planning with Latent. Language]

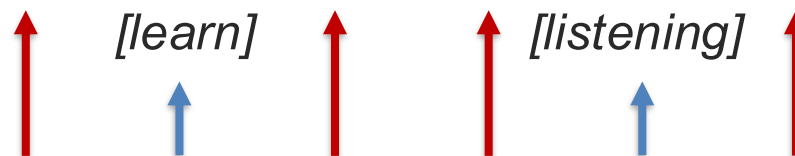
Co-learning for Pre-training

Predicting images from corresponding language

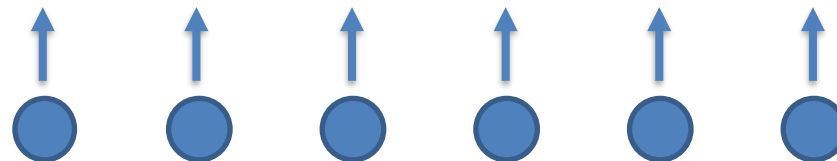
Voken (visual token) classification



Masked language modeling



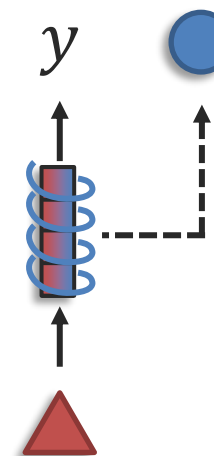
BERT language model



Humans [mask] language by [mask] speaking

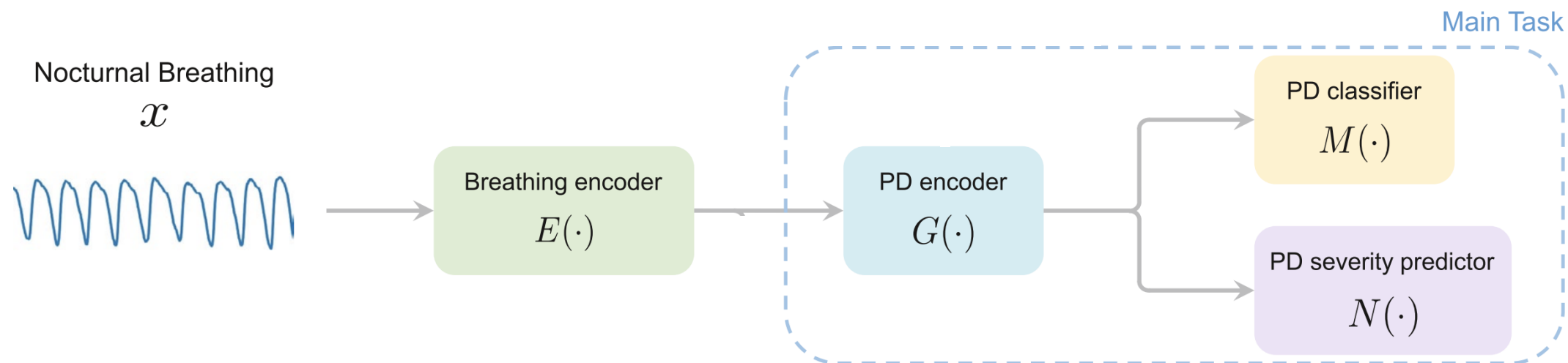
Only text used at test-time

Multimodal co-learning > language-only training



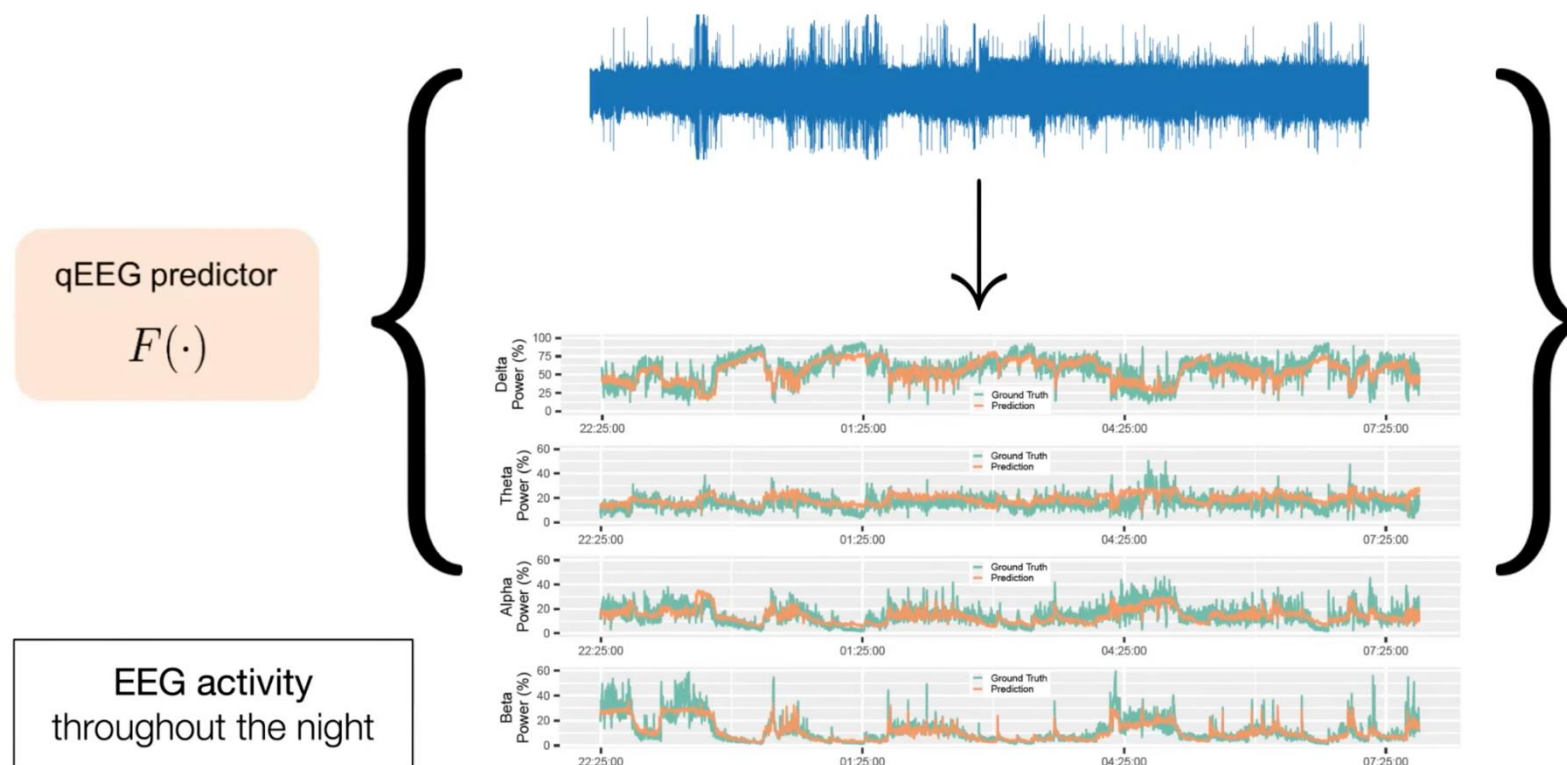
Co-learning for Dense Supervision

10hours of breathing data to detect Parkinson's Disease (PD) – sparse 1 bit signal



Co-learning for Dense Supervision

Predicting paired EEG data as auxiliary task provides dense supervision



Limits of Co-learning

Vision-language pretrained models on lexical grounding

Sentence-level semantic tasks

Encoder	SRL	Coref.	SPR	Rel.
BERT _{base}	90.10 ± 0.20	95.90 ± 0.00	83.70 ± 0.00	76.25 ± 0.05
VideoBERT _{text}	84.33 ± 0.05	92.47 ± 0.05	78.23 ± 0.05	65.83 ± 0.21
VideoBERT _{VL}	84.73 ± 0.05	92.82 ± 0.05	78.80 ± 0.00	66.37 ± 0.80
VisualBERT _{text}	89.00 ± 0.00	94.87 ± 0.05	82.27 ± 0.05	74.37 ± 0.19
VisualBERT _{VL}	89.57 ± 0.21	95.13 ± 0.05	82.17 ± 0.09	74.83 ± 0.05

Not much improvements with visual co-learning

Semantic Role Labeling “The **carrots** are then pureed in the food processor”
 Entity Coreference “After the **apples** are chopped, put **them** in the bowl”

Limits of Co-learning

Vision-language pretrained models on seemingly multimodal tasks

Physical commonsense QA

Encoder	Linear	MLP	Trans.
BERT _{base}	55.43 ± 0.31	57.98 ± 0.16	60.12 ± 1.43
VideoBERT _{text}	57.87 ± 0.64	58.97 ± 0.44	62.35 ± 1.23
VideoBERT _{VL}	58.51 ± 0.20	58.56 ± 0.27	63.66 ± 1.31
VisualBERT _{text}	54.81 ± 0.19	56.81 ± 0.24	58.63 ± 0.79
VisualBERT _{VL}	55.83 ± 0.27	59.10 ± 0.11	61.66 ± 1.08

Marginal improvements with visual co-learning

“How to remove gloss from furniture?”

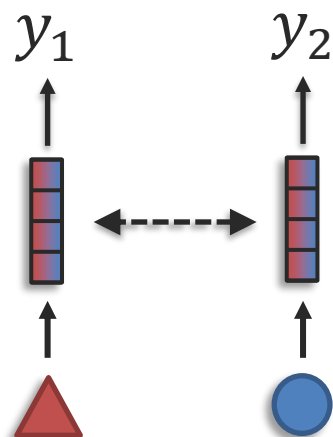


“Rub furniture with steel wool/cotton ball”

Part 3: Model Induction

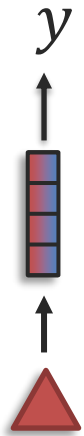
Definition: Keeping individual unimodal models separate but inducing common behavior across separate models.

Model Induction



Self-training

Warmup: a single view – Self-training



Assume:

1. Labeled data $\{X_1^L, Y\}$.
2. Unlabeled data $\{X_1^U\}$.

Train:

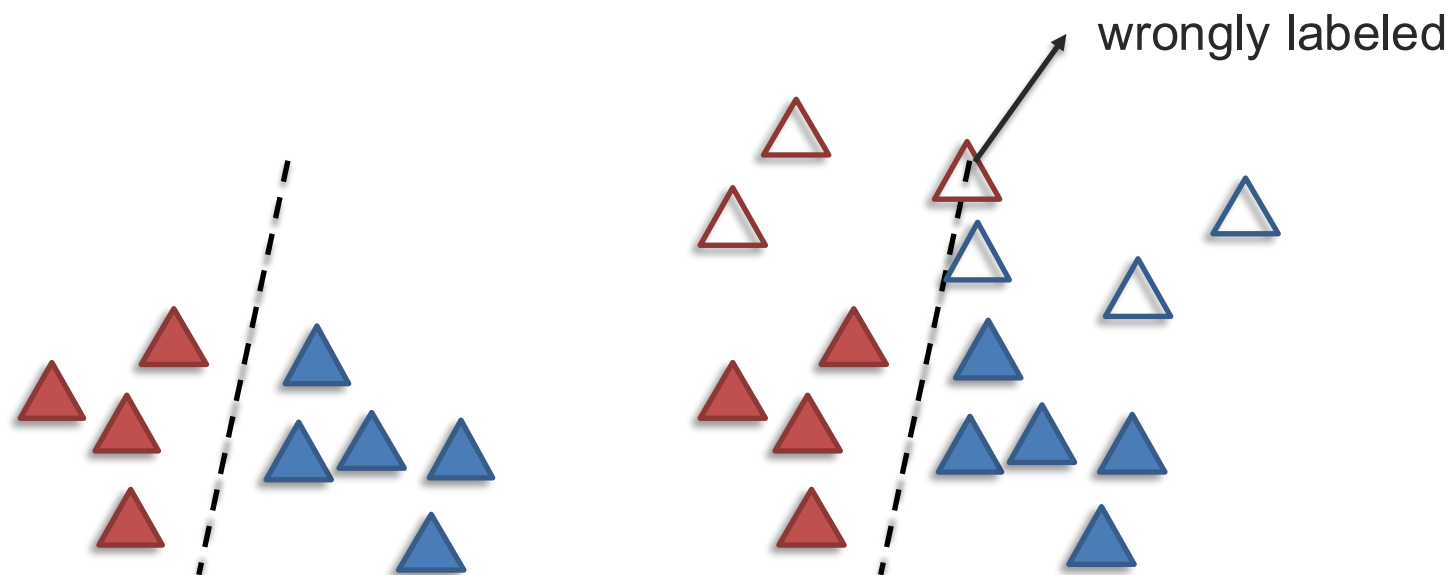
1. Train classifier f_1 on $\{X_1^L, Y\}$.
2. Use classifier f_1 to label the most confident examples in $\{X_1^U\}$ and add it to the labeled set $\{X_1^U, Y = f_1(X_1^U)\}$.
3. Go to 1, and repeat until there are no more unlabeled samples.

Test:

1. For a new unlabeled sample $\{X_1\}$, output $f_1(X_1)$.

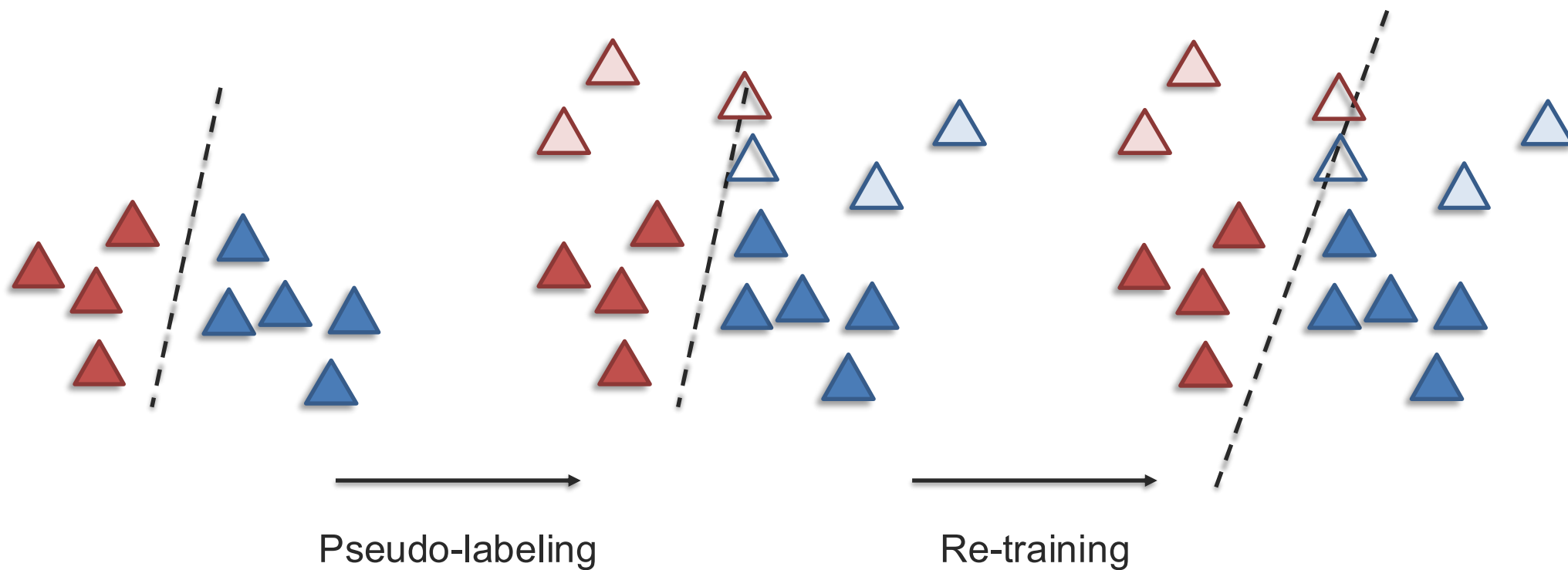
Self-training

Warmup: a single view – Self-training



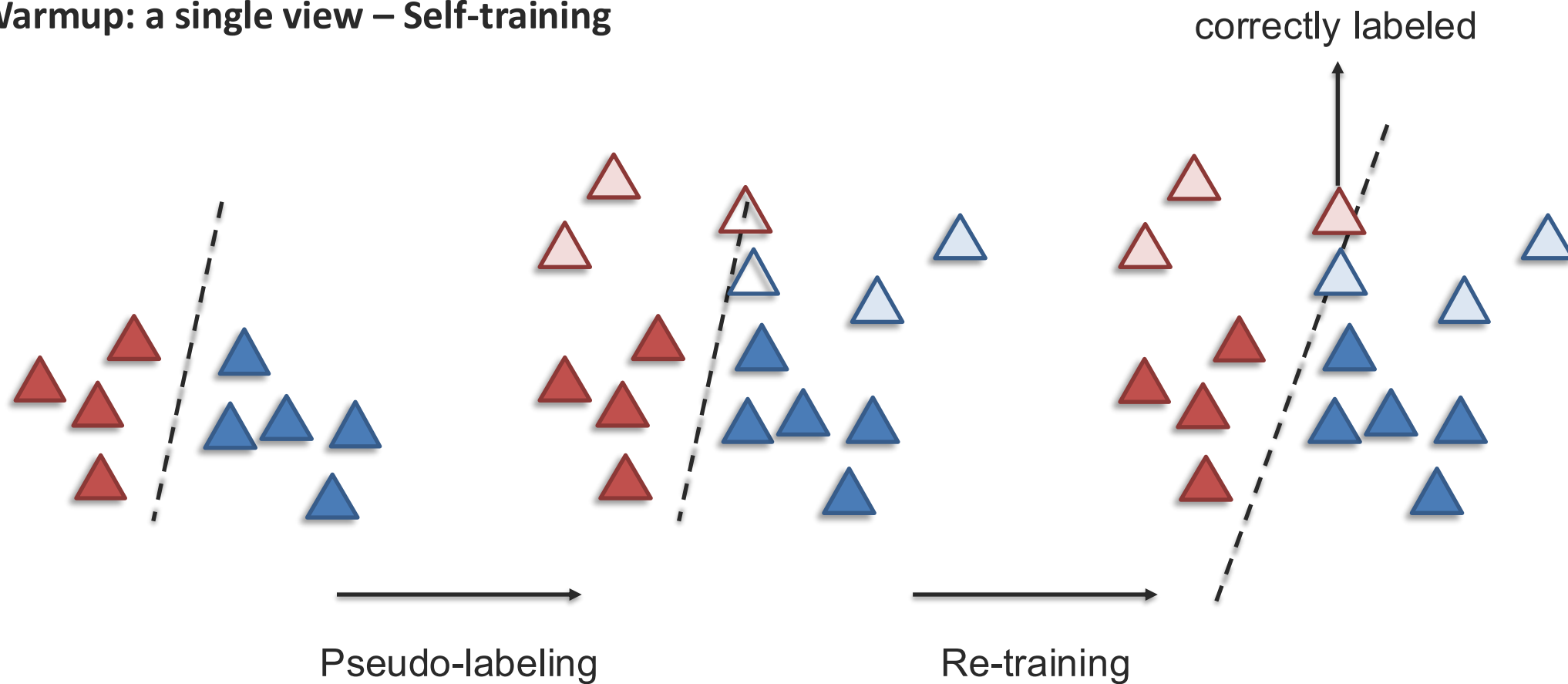
Self-training

Warmup: a single view – Self-training



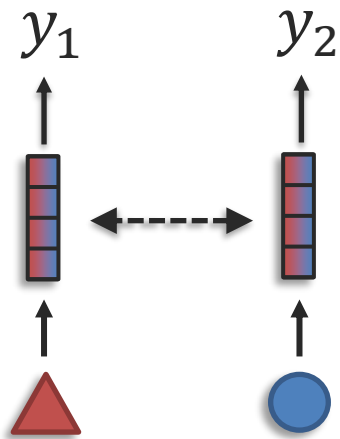
Self-training

Warmup: a single view – Self-training



Self-training

From self-training to co-training



Ingredients:

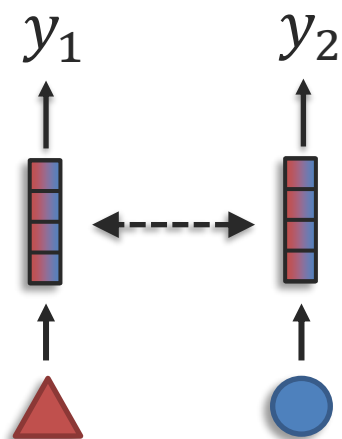
- ↯ Two views on the data: x_1 and x_2
- ↯ Two classifiers: $x_1 \rightarrow y$ and $x_2 \rightarrow y$
- ↯ A bit of labeled data (x_1, x_2, y) ; lots of unlabeled data (x_1, x_2)

Assumptions:

1. Multi-view redundancy: either view is sufficient to predict the label alone, with enough data.

Co-training

Algorithm



Assume:

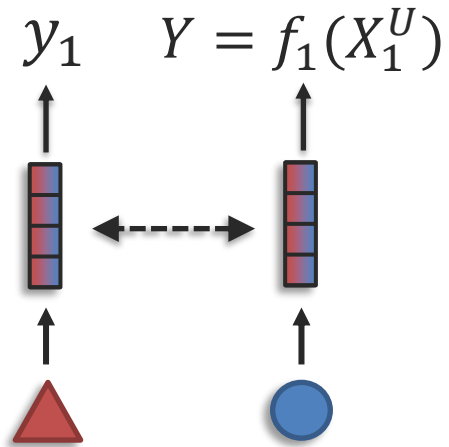
1. **Small** amount of labeled data $\{X_1^L, X_2^L, Y\}$.
2. **Lots** of unlabeled data $\{X_1^U, X_2^U\}$.

Train:

1. Train classifier f_1 on $\{X_1^L, Y\}$ and f_2 on $\{X_2^L, Y\}$.

Co-training

Algorithm



Assume:

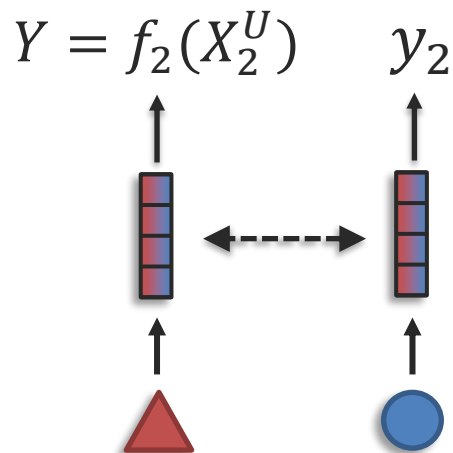
1. **Small** amount of labeled data $\{X_1^L, X_2^L, Y\}$.
2. **Lots** of unlabeled data $\{X_1^U, X_2^U\}$.

Train:

1. Train classifier f_1 on $\{X_1^L, Y\}$ and f_2 on $\{X_2^L, Y\}$.
2. Use classifier f_1 to label the most confident examples in $\{X_1^U\}$ and add it to the labeled set to train f_2 $\{X_2^U, Y = f_1(X_1^U)\}$.

Co-training

Algorithm



Assume:

1. **Small** amount of labeled data $\{X_1^L, X_2^L, Y\}$.
2. **Lots** of unlabeled data $\{X_1^U, X_2^U\}$.

Train:

1. Train classifier f_1 on $\{X_1^L, Y\}$ and f_2 on $\{X_2^L, Y\}$.
2. Use classifier f_1 to label the most confident examples in $\{X_1^U\}$ and add it to the labeled set to train f_2 $\{X_2^U, Y = f_1(X_1^U)\}$.
3. Use classifier f_2 to label the most confident examples in $\{X_2^U\}$ and add it to the labeled set to train f_1 $\{X_1^U, Y = f_2(X_2^U)\}$.
4. Go to 1, and repeat until there are no more unlabeled samples.

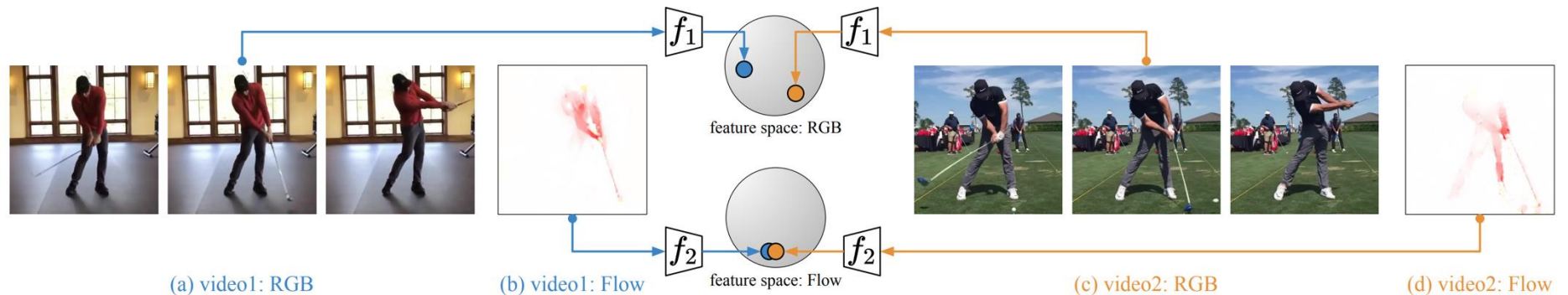
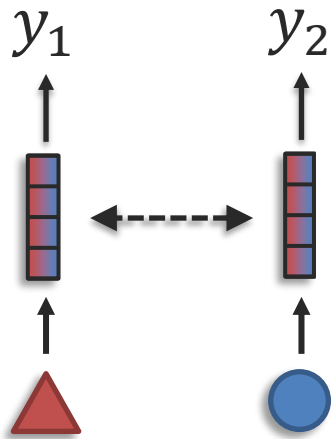
Test:

1. For a new unlabeled sample $\{X_1, X_2\}$, ensemble $f_1(X_1)$ and $f_2(X_2)$.

Modern Co-training

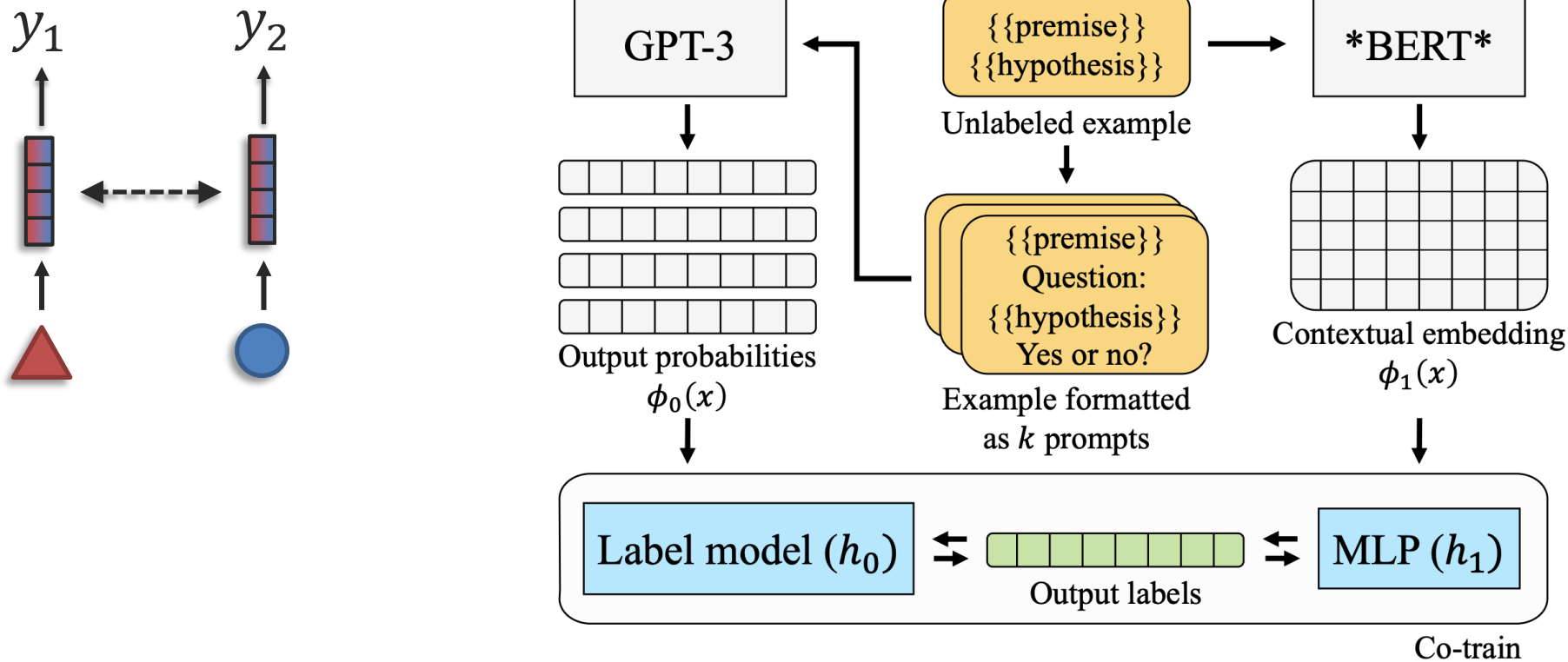
Co-training between RGB and optical flow for activity recognition

- Positive samples hard to discover in RGB space can be easily found in flow space, and vice-versa (e.g., RGB sensitive to background differences but not flow).
- Can use co-training between 2 RGB and flow contrastive learning modules.



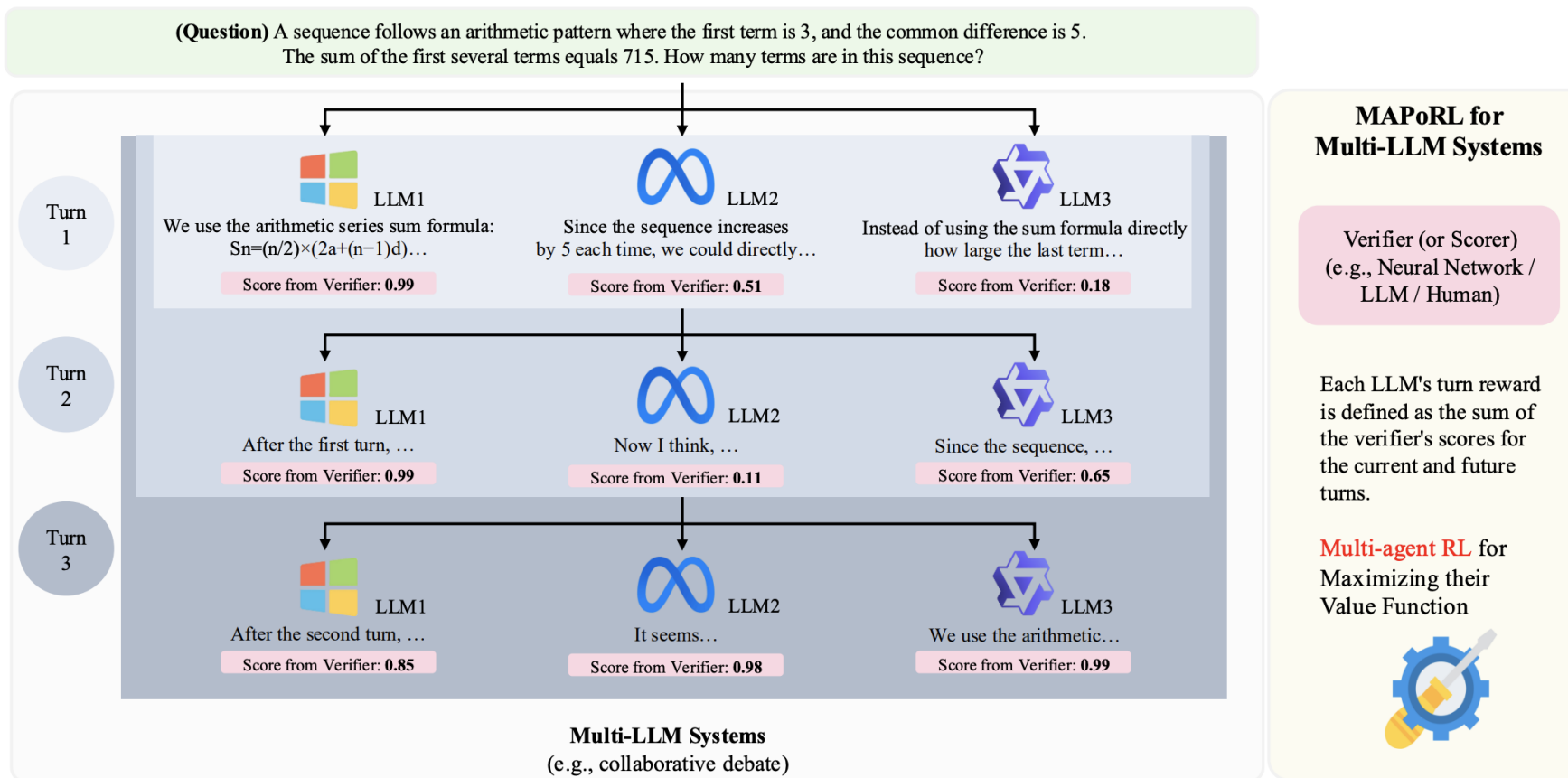
Modern Co-training

Language-model prompting



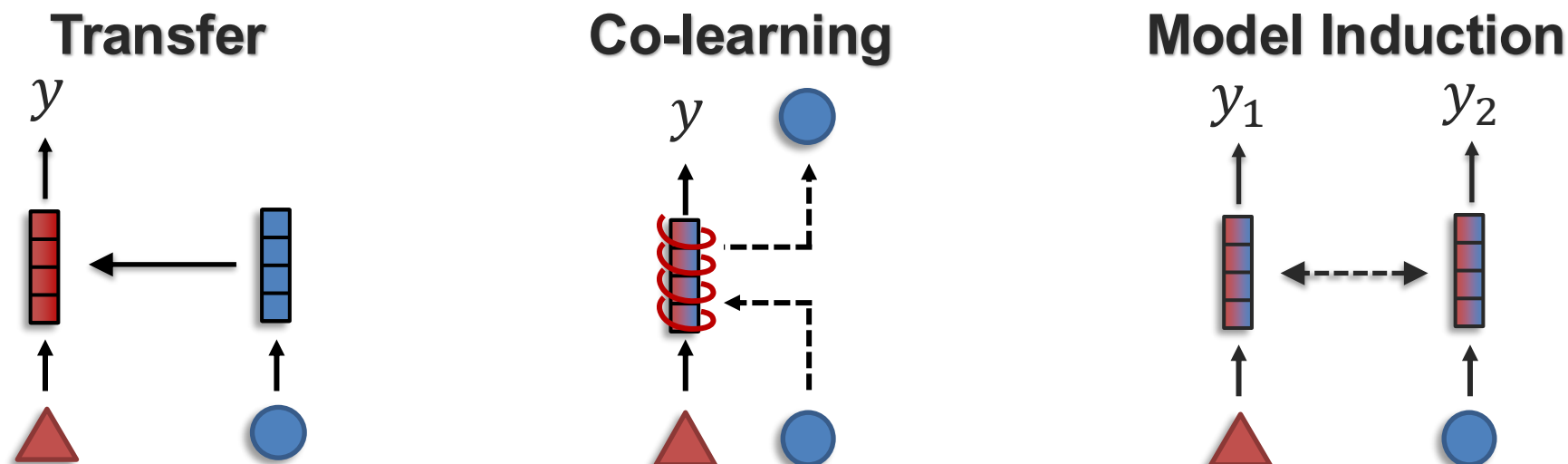
Modern Co-training

Multi-agent LLMs, debate, co-training



Summary: How to Cross-modal Transfer

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



1. Decide on secondary modalities
2. Decide on auxiliary input or auxiliary output
3. Decide on modifying model or using APIs only

Today's lecture

- 1 Basics of cross-modal transfer
- 2 Cross-modal transfer via fusion
- 3 Cross-modal transfer via alignment
- 4 Cross-modal transfer via translation