**Language Technologies Institute**

# Multimodal Machine Learning

## Lecture 11.2: Transference 2 – Co-learning and Co-training

**Paul Liang**

# Reading Assignments are Back!

- Four main steps for the reading assignments
  - Monday 8pm: Official start of the assignment
  - Wednesday 8pm: Select your paper
  - **Friday 8pm:** Post your summary
  - **Monday 8pm:** Post your extra comments (5 posts)
- **4 papers:** multimodal multi-hop reasoning, multimodal geometric reasoning, multimodal robotics, multimodal knowledge bases.

# Final Project Report (Due Sunday 12/11 at 8pm)

Main goals:

1. Produce a research paper which will motivate your research problem, describe the prior work, present your research contributions, explain the details of your experiments, and discuss your results.

2. Novel research ideas (N-1 new ideas for N students)

   - Novel algorithm
   - Novel application

3. Incorporate feedback from previous milestones

4. Compare to multimodal baselines from midterm report

   1. Did the proposed ideas solve the errors highlighted in error analysis?
   2. Broader implications of proposed ideas.

# Final Project Presentations (Tuesday 12/6 and Thursday 12/8)

Main objective:

- Present your research ideas and get feedback from classmates
- Focus on only one of your new research ideas
- All students should present and answer questions
- Be sure to be on time! We have many presentations each day ☺
- All presentations are in person (no remote presentations)

Presentation length:

- 30-seconds elevator pitch
- 4-minute full presentation – all students should present

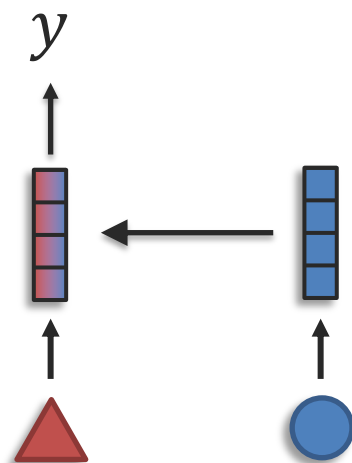- Following each presentation, audience will be asked to share feedback
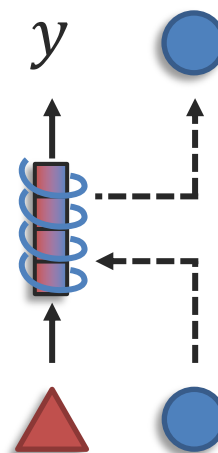
# Transference

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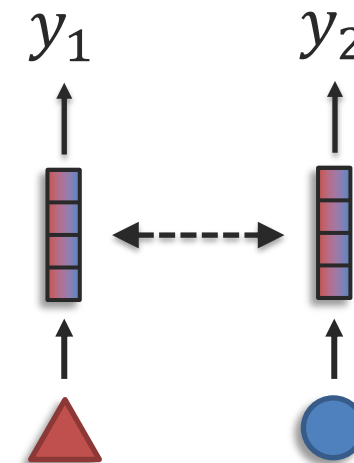
**Sub-challenges:**
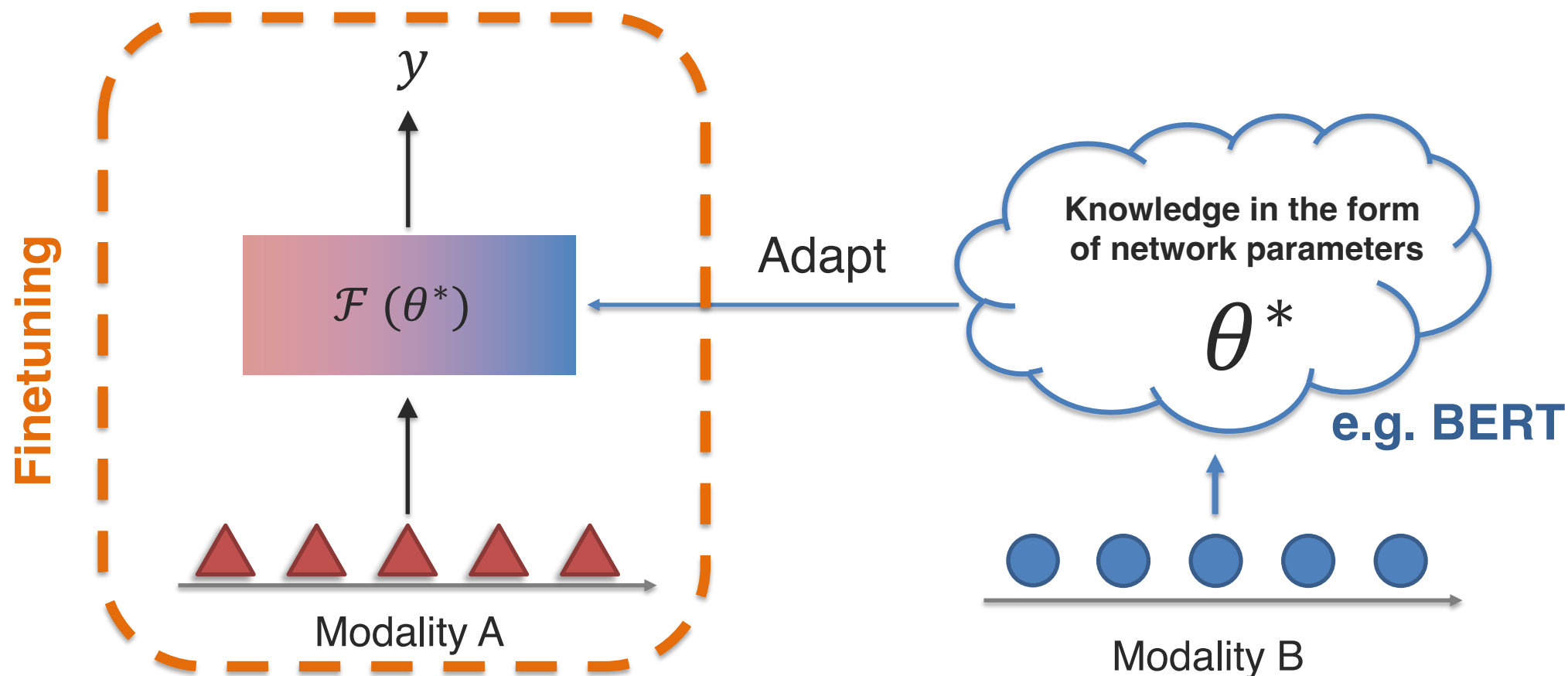


**Transfer**         **Co-learning**         **Model Induction**
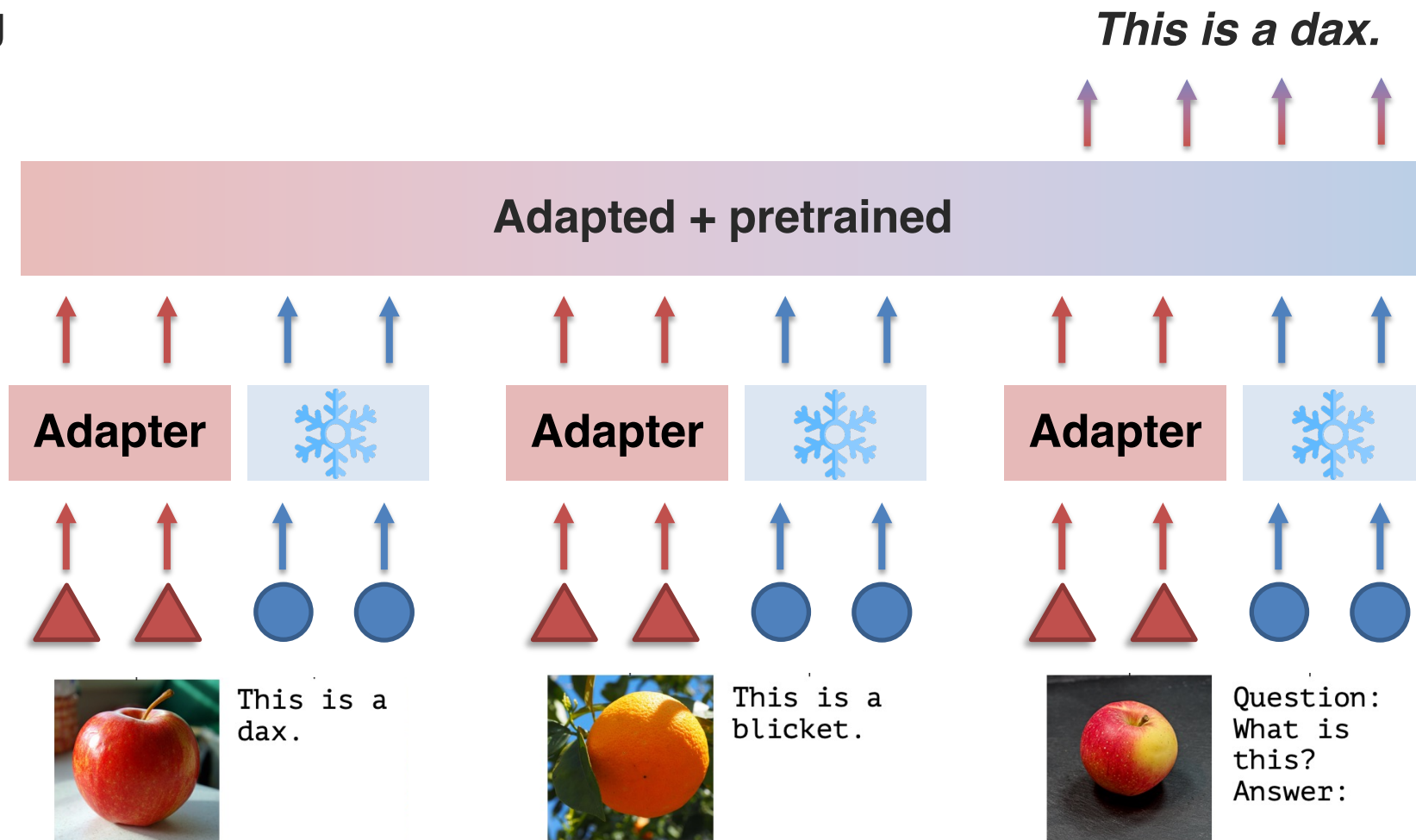
**Definition:** Transferring knowledge from large-scale pretrained models to downstream tasks involving the primary modality.

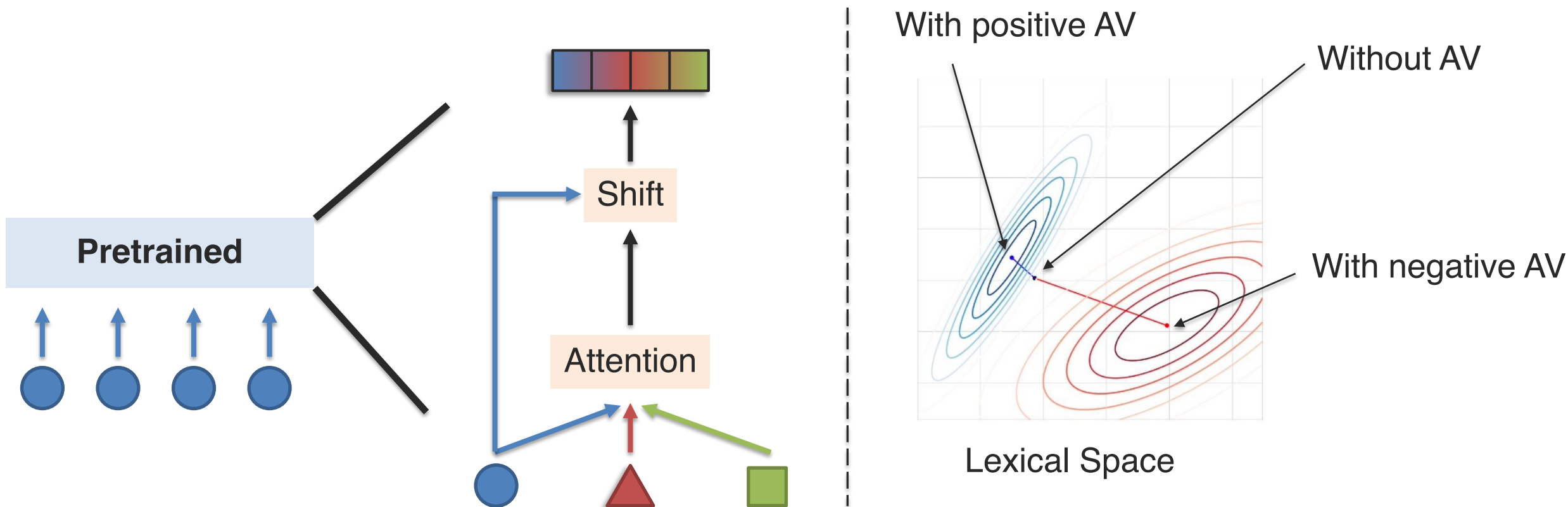# Sub-Challenge 5a: Transfer via Pretrained Models

**Transfer via prefix tuning**

*This is a dax.*

Few-shot image classification:



[Tsimpoukelli et al., Multimodal Few-Shot Learning with Frozen Language Models. NeurIPS 2021]

**Transfer via representation tuning**



With positive AV

Without AV

With negative AV

Lexical Space

Shift

Attention

Pretrained

[Ziegler et al., Encoder-Agnostic Adaptation for Conditional Language Generation. arXiv 2019]

[Rahman et al., Integrating Multimodal Information in Large Pretrained Transformers. ACL 2020]

# Sub-Challenge 5a: Transfer via Pretrained Models

1. Disentanglement

$$\mathcal{L}_\beta(x) = \mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x|z)] - \beta \cdot \mathrm{KL}(q_\phi(z|x)||p(z))$$

2. Conditioning

$$p(\boldsymbol{x}_{0:T} \mid y) = p(\boldsymbol{x}_T) \prod_{t=1}^{T} p_{\boldsymbol{\theta}}(\boldsymbol{x}_{t-1} \mid \boldsymbol{x}_t, y)$$

3. Prompt tuning

4. Representation tuning

5. Classifier gradient tuning

$$\nabla \log p(\boldsymbol{x}_t \mid y) = \underbrace{\nabla \log p(\boldsymbol{x}_t)}_{\text{unconditional score}} + \gamma \underbrace{\nabla \log p(y \mid \boldsymbol{x}_t)}_{\text{classifier gradient}}$$

6. Classifier-free tuning

$$\nabla \log p(\boldsymbol{x}_t \mid y) = \gamma \underbrace{\nabla \log p(\boldsymbol{x}_t \mid y)}_{\text{conditional score}} + (1 - \gamma) \underbrace{\nabla \log p(\boldsymbol{x}_t)}_{\text{unconditional score}}$$

# Multitask and Transfer Learning

**How can we transfer knowledge across multiple tasks, each over a different subset of modalities?**



Video classification

Language    Video    Audio

Sentiment, emotions

Audio    Video

Robot dynamics

Video    Time-series

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

[Liang et al., HighMMT: Towards Modality and Task Generalization for High-Modality Representation Learning. arXiv 2022]

# Multitask and Transfer Learning

**Transfer across partially observable modalities**

HighMMT: unified model + parameter sharing + multitask and transfer learning



**Non-parallel multitask learning**

**Task-specific classifiers**

**Same model architecture!**

**Shared multimodal model**

**Same parameters!**

**Modality-specific embeddings**

**Standardized input sequence**

[Liang et al., HighMMT: Towards Modality and Task Generalization for High-Modality Representation Learning. arXiv 2022]

# Multitask and Transfer Learning

**Traditional approaches: different model + different parameters**



Image-text retrieval   Design interface   Robotic manipulation   Disease codes   Emotions   Sarcasm   Humor

Language   Image   Audio   Video   Sensors   Proprioception   Speech   Time-series   Set   Table

Performance →   Efficiency (params) →

● All model combinations (>10,000)
● Pareto front

[Liang et al., HighMMT: Towards Modality and Task Generalization for High-Modality Representation Learning. arXiv 2022]

# Multitask and Transfer Learning

**Traditional approaches: different model + different parameters**



[Liang et al., HighMMT: Towards Modality and Task Generalization for High-Modality Representation Learning. arXiv 2022]

# Multitask and Transfer Learning

**Traditional approaches: different model + different parameters**

Image-text retrieval    Design interface    Robotic manipulation    Disease codes    Emotions    Sarcasm    Humor

## HighMMT multitask model

Language    Image    Audio    Video    Sensors    Proprioception    Speech    Time-series    Set    Table

- ● All model combinations (>10,000)
- ● Pareto front
- ○ HighMMT single-task
- ● HighMMT multitask

[Liang et al., HighMMT: Towards Modality and Task Generalization for High-Modality Representation Learning. arXiv 2022]

# Multitask and Transfer Learning

**HighMMT heterogeneity-aware sharing: estimate heterogeneity to determine parameter sharing**



Image-text retrieval · Design interface · Robotic manipulation · Disease codes · Emotions · Sarcasm · Humor

**HighMMT heterogeneity-aware sharing**

Language · Image · Audio · Video · Sensors · Proprioception · Speech · Time-series · Set · Table

- All model combinations (>10,000)
- Pareto front
- HighMMT single-task
- HighMMT multitask
- **HighMMT heterogeneity-aware**

[Liang et al., HighMMT: Towards Modality and Task Generalization for High-Modality Representation Learning. arXiv 2022]

# Multitask and Transfer Learning

**Transfer across partially observable modalities**
HighMMT: unified model + parameter sharing + multitask and transfer learning



[Liang et al., HighMMT: Towards Modality and Task Generalization for High-Modality Representation Learning. arXiv 2022]

# Multitask and Transfer Learning

**Transfer across partially observable modalities**
HighMMT: unified model + parameter sharing + multitask and transfer learning

Target task: MIMIC

67.7%    68.3%    68.5%    **68.5%**

# source tasks    0    1    2    3

(from different modalities, research
areas, and tasks)

Target task: UR-FUNNY

63.3%    64.1%    65.5%    **65.7%**

# source tasks    0    1    2    3

(from different modalities, research
areas, and tasks)

Achieves both multitask and transfer capabilities across modalities and tasks

[Liang et al., HighMMT: Towards Modality and Task Generalization for High-Modality Representation Learning. arXiv 2022]

# Multitask and Transfer Learning

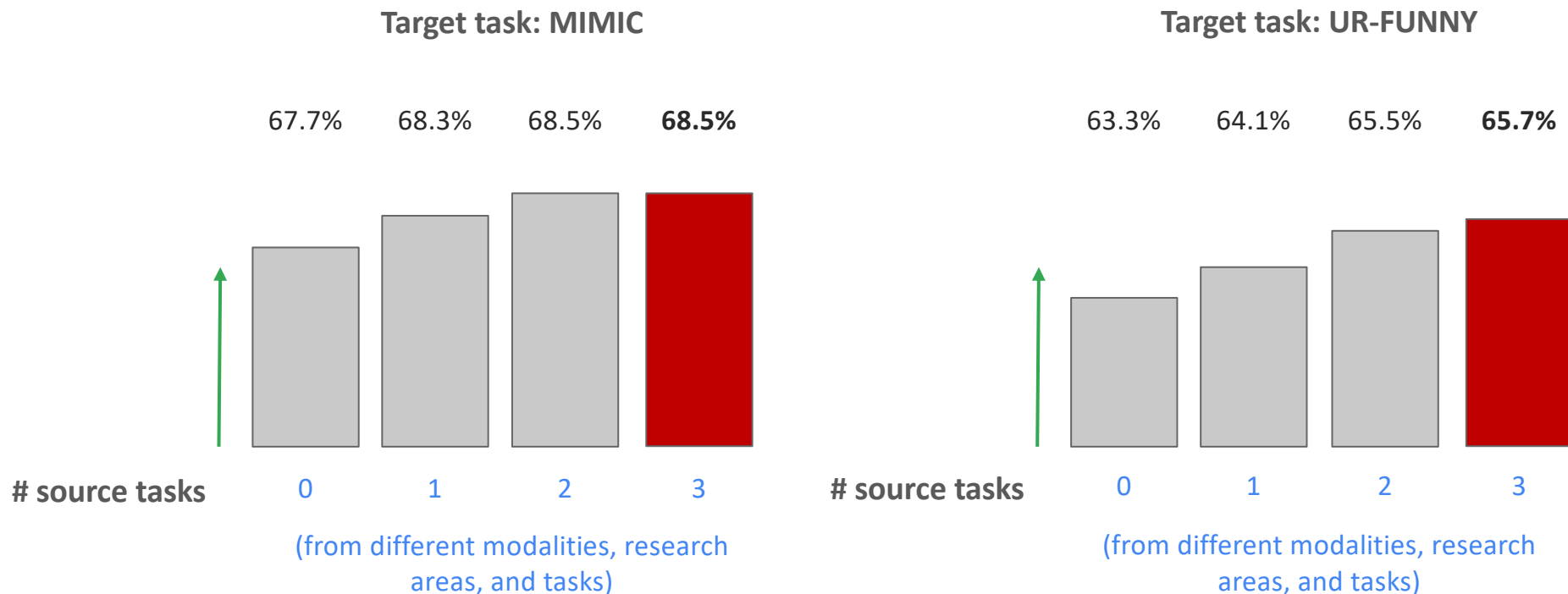**Transfer across partially observable modalities**
Gato: unified model + parameter sharing + multitask learning



Language modeling     Visual QA     Atari games     Robot manipulation

**Non-parallel multitask learning**

**Task-specific classifiers**

**Same model architecture!**

**Shared multimodal model**

**Same parameters!**

**Modality-specific embeddings**

**Standardized input sequence**

I'm going to London

Q: What's in the picture?
A: It's a cute cat

Text     Image     Text     Image     Action     Image     Proprioception     Action

[Reed et al., A Generalist Agent. arXiv 2022]

# Multitask and Transfer Learning

**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.



Video classification    Sentiment, emotions    Robot dynamics

**HighMMT/Gato**

**Shared multimodal model?**

**Modality-specific embeddings?**

**Standardized input sequence?**

Language    Video    Audio        Audio    Video        Video    Time-series
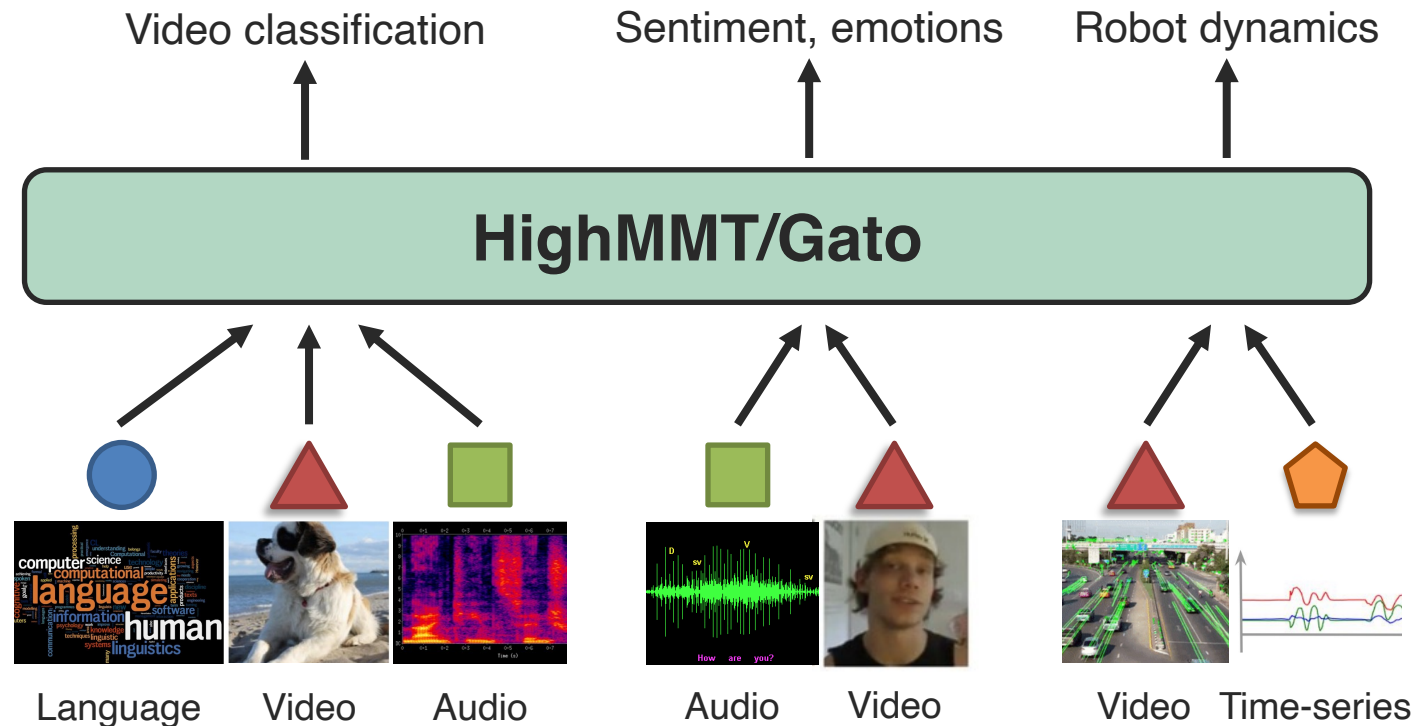
[Liang et al., HighMMT: Towards Modality and Task Generalization for High-Modality Representation Learning. arXiv 2022]

# Multitask and Transfer Learning

**Many more dimensions of transfer**



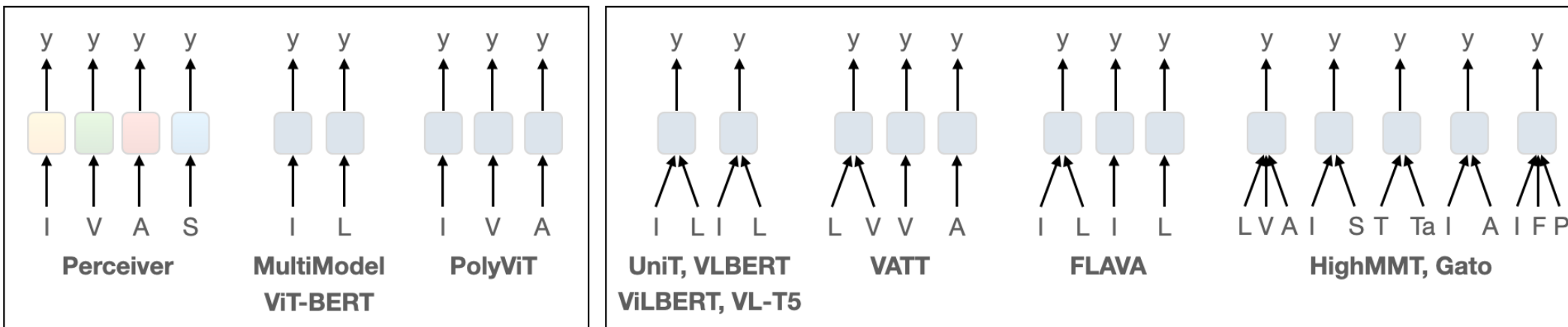**Unified encoder for unimodal learning**

Perceiver — I V A S

MultiModel — I L
ViT-BERT

PolyViT — I V A

**Multimodal multitask learning**

UniT, VLBERT ViLBERT, VL-T5 — I L I L

VATT — L V V A

FLAVA — I L I L

HighMMT, Gato — L V A I  S T  Ta I  A I F P

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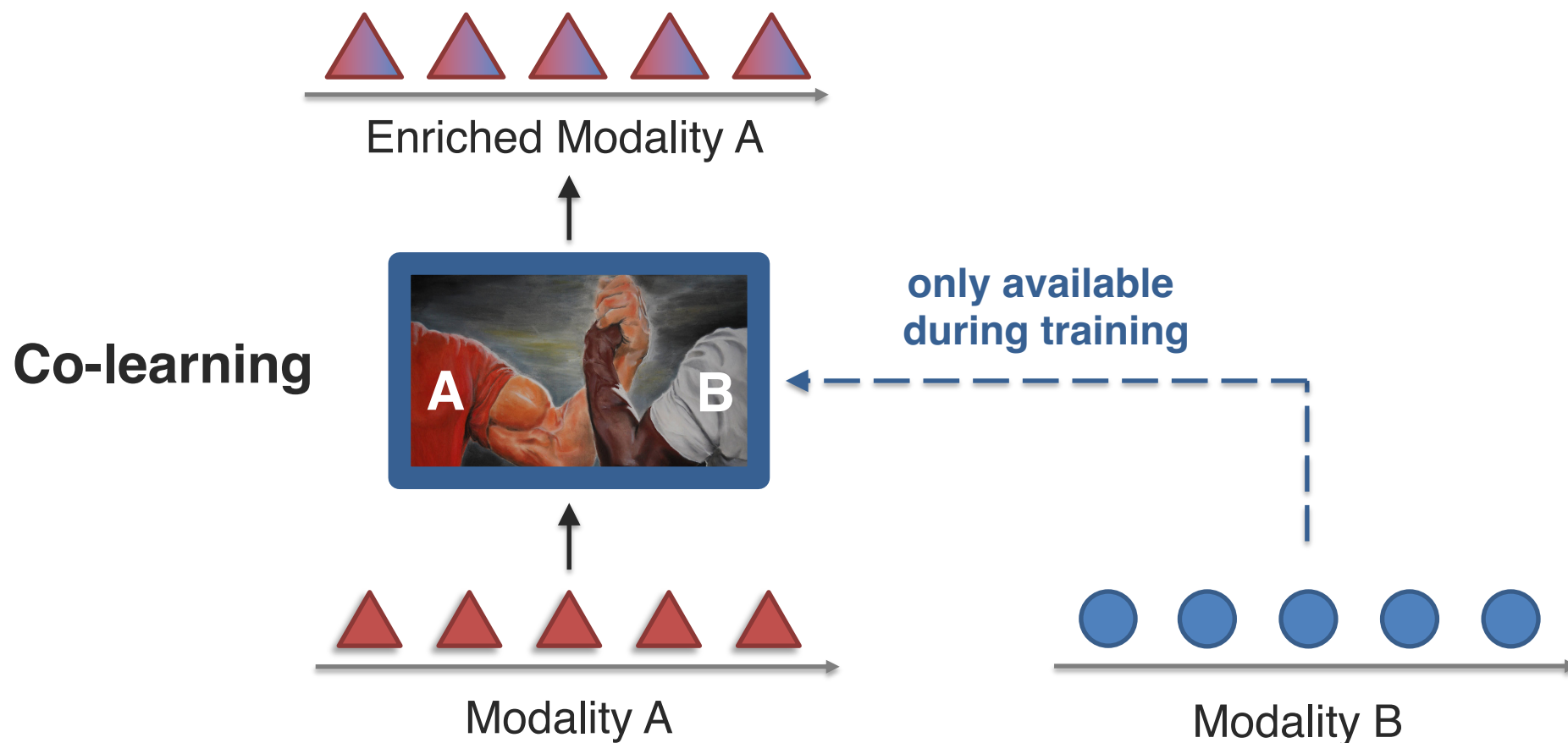
common architecture

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
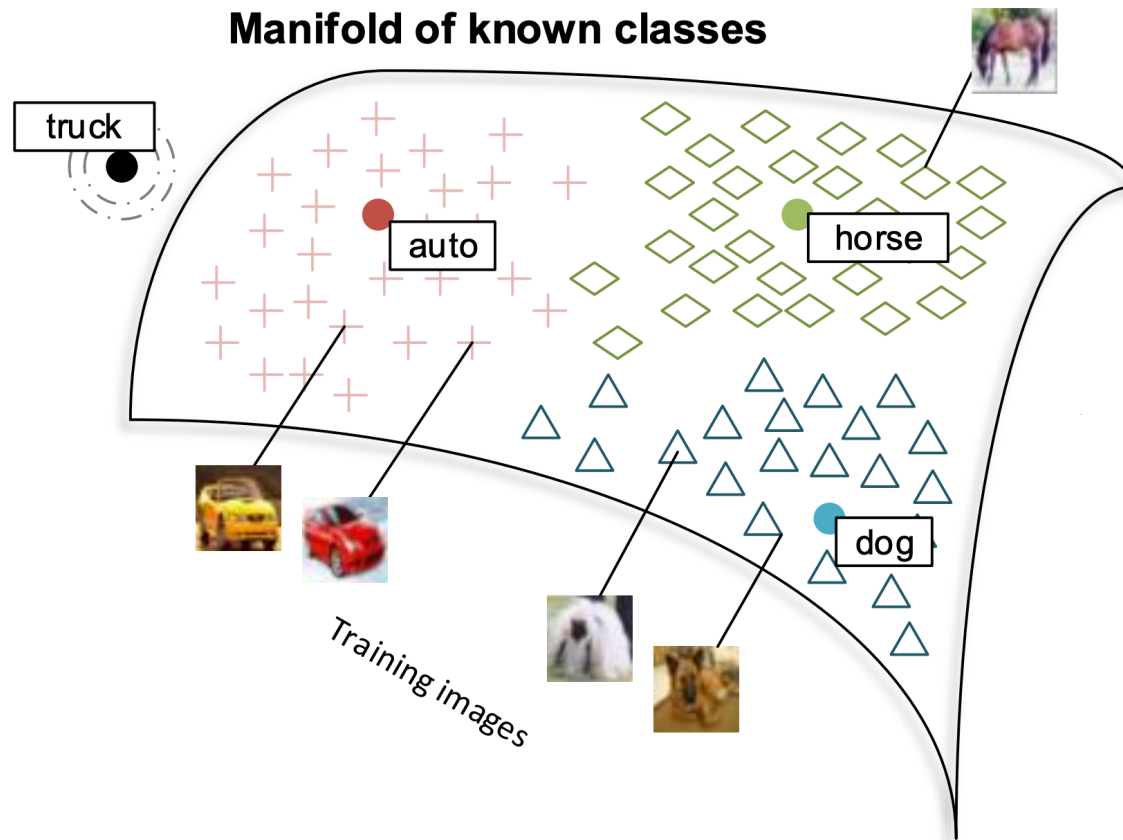- Complexity in data, modeling, and training
- Interpretability (next section)

# Sub-Challenge 5b: Co-learning

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



Enriched Modality A

**Co-learning**

A          B

**only available during training**

Modality A

Modality B

# Co-learning via Representation

**Representation coordination: word embedding space for zero-shot visual classification**



**Manifold of known classes**

truck

auto

horse

dog

Training images

Recall representation coordination!

encoder $f_A$

encoder $f_B$

$g(\mathbf{z}_A, \mathbf{z}_B)$

$\mathbf{z}_B$

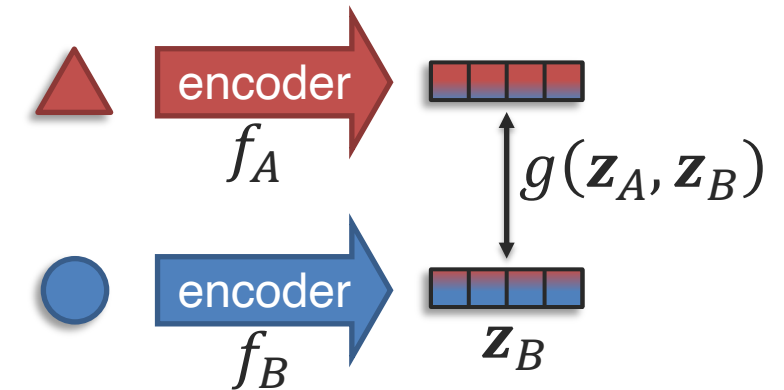[Socher et al., Zero-Shot Learning Through Cross-Modal Transfer. NeurIPS 2013]

# Co-learning via Representation

**Representation coordination: word embedding space for zero-shot visual classification**



**Manifold of known classes**

truck

auto

horse

New test image from unknown class

dog

cat

Training images

Recall representation coordination!

encoder $f_A$

encoder $f_B$

$g(\mathbf{z}_A, \mathbf{z}_B)$

$\mathbf{z}_B$

Only images used at test-time
Enables zero-shot image classification
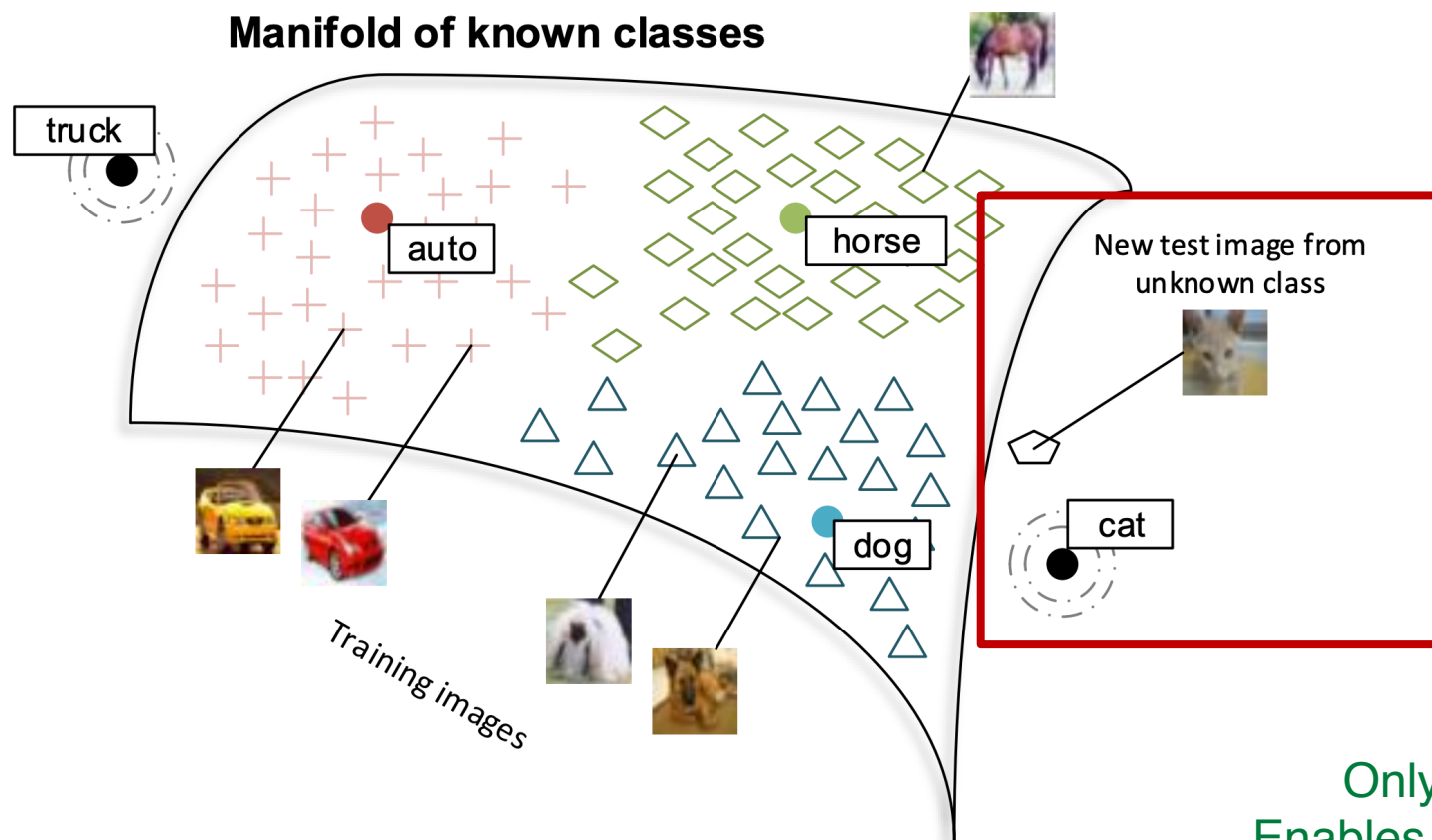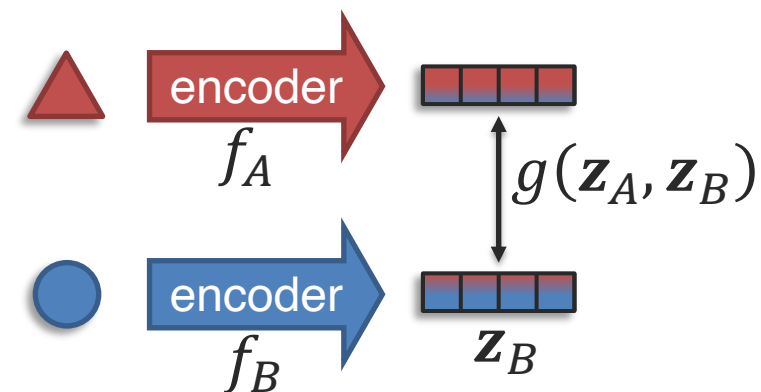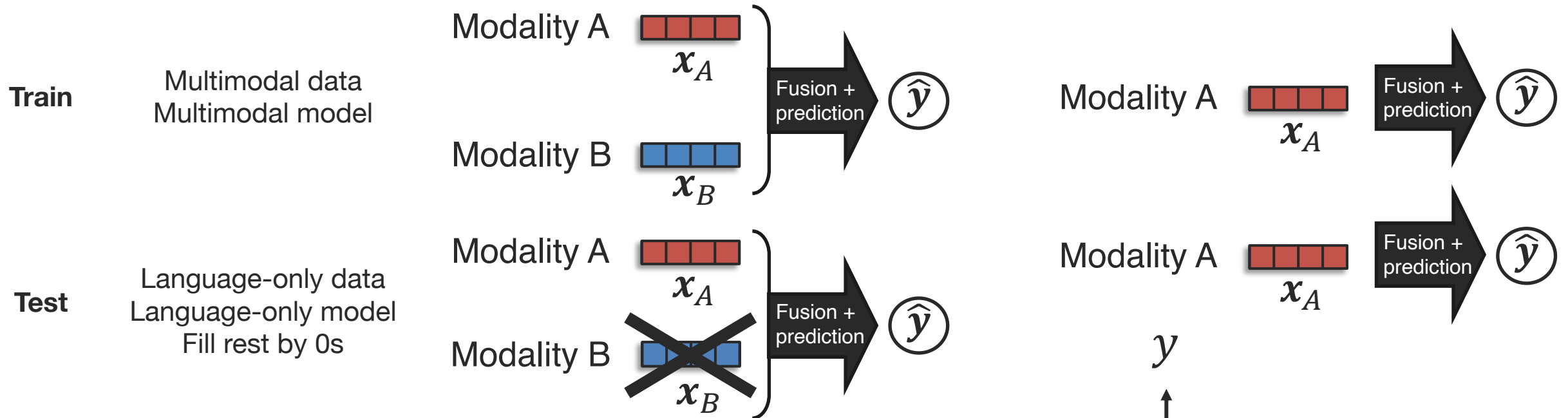
[Socher et al., Zero-Shot Learning Through Cross-Modal Transfer. NeurIPS 2013]

# Co-learning via Representation

**Representation fusion**    **Multimodal co-learning**    **Unimodal learning**



Only text used at test-time

Multimodal co-learning > language-only training

[Zadeh et al., Foundations of Multimodal Co-learning. Information Fusion 2020]

# Co-learning via Generation

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



Enriched Modality A

Modality B

**Co-learning**

A    B

**only available during training**

Modality A

# Co-learning via Generation
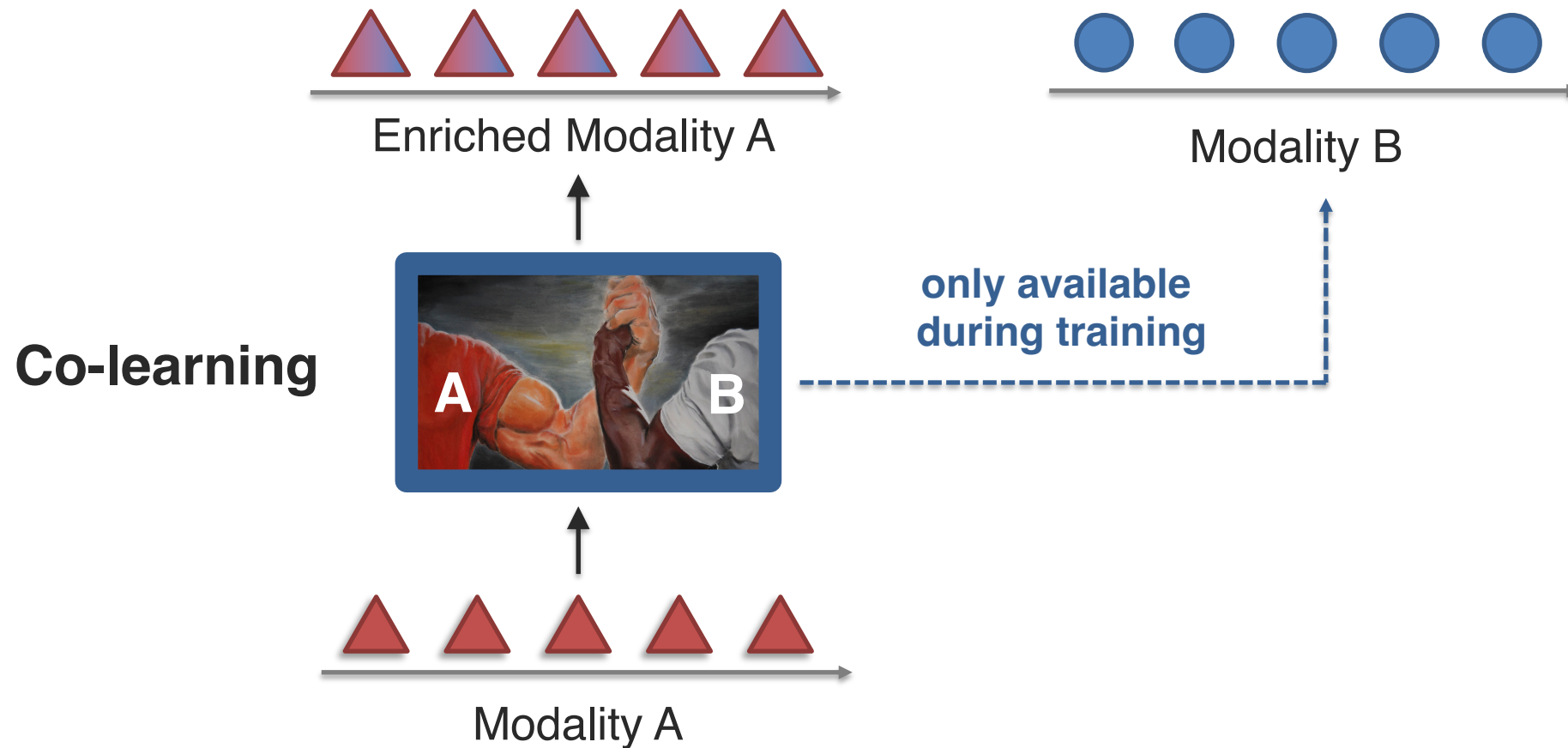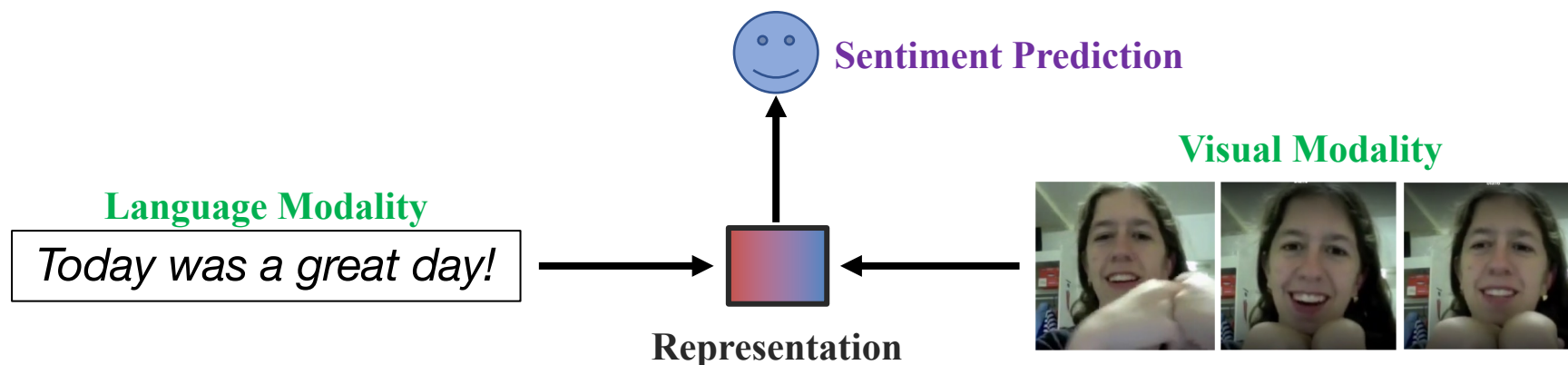
**Bimodal translations**



**Sentiment Prediction**

**Language Modality**

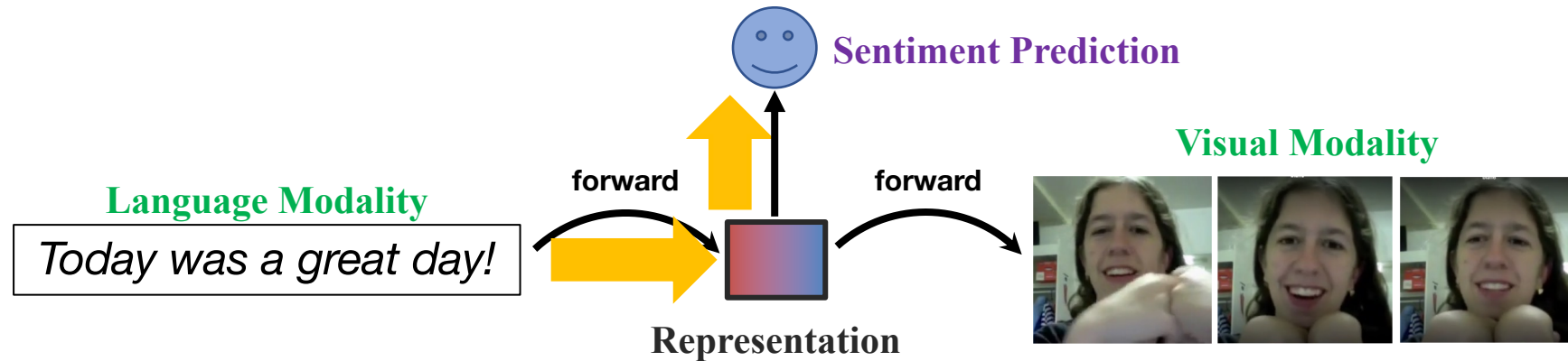*Today was a great day!*

**Representation**

**Visual Modality**

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.

[Pham et al., Found in Translation: Learning Robust Joint Representations via Cyclic Translations Between Modalities. AAAI 2019]

# Co-learning via Generation

**Bimodal translations**



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

[Pham et al., Found in Translation: Learning Robust Joint Representations via Cyclic Translations Between Modalities. AAAI 2019]

# Co-learning via Generation

**Bimodal translations**



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

[Pham et al., Found in Translation: Learning Robust Joint Representations via Cyclic Translations Between Modalities. AAAI 2019]
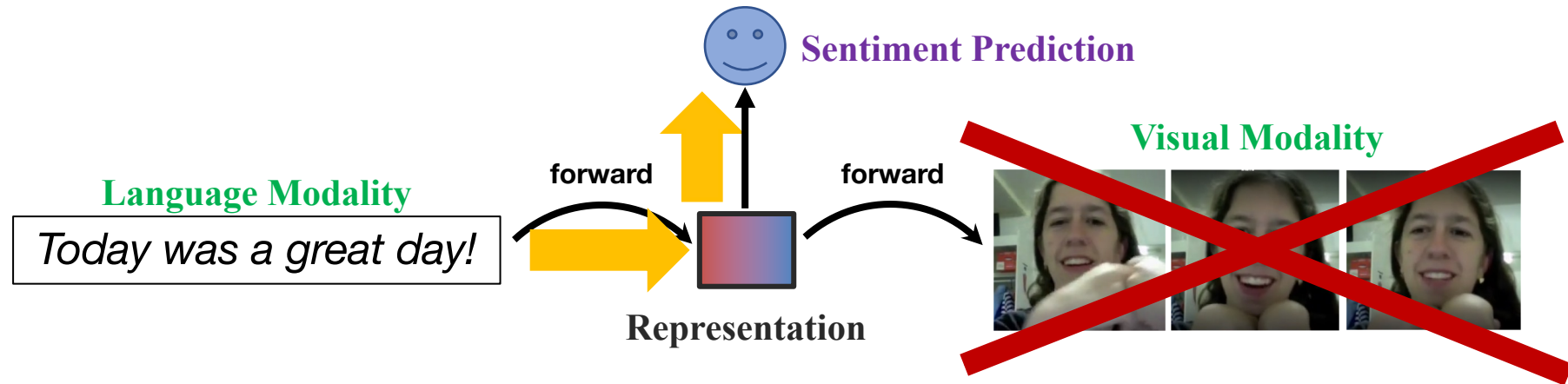
# Co-learning via Generation

**Bimodal cyclic translations**



**Sentiment Prediction**

**Language Modality**

*Today was a great day!*

forward  forward

backward  backward

**Visual Modality**

Solution: cyclic translations from visual back to language

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

[Pham et al., Found in Translation: Learning Robust Joint Representations via Cyclic Translations Between Modalities. AAAI 2019]

# Co-learning via Generation

**Predicting images from corresponding language**

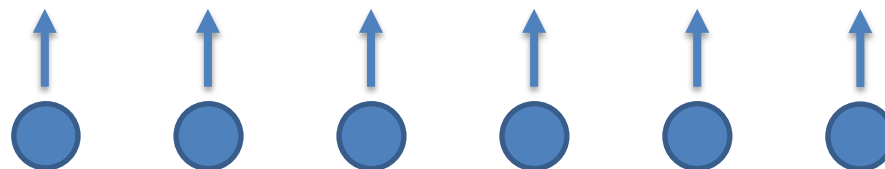**Voken (visual token) classification**

**Masked language modeling**

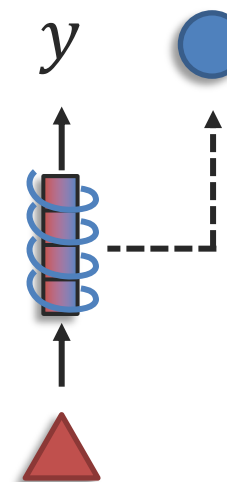*[learn]*          *[listening]*

**BERT language model**

$y$

*Humans [mask] language by [mask] speaking*

Only text used at test-time

Multimodal co-learning > language-only training

[Tan and Bansal, Vokenization: Improving Language Understanding with Contextualized, Visual-Grounded Supervision. EMNLP 2020]

# Co-learning via Generation

**Language** ▲

**Visual**
(image) ●

**Information primarily in language modality**

- Syntactic structure
- Vocabulary, morphology
- …

**Information in both modalities**

- Described people, objects, actions
- Illustrative gestures, motion
- …

**Information primarily in visual modality**

- Texture, visual appearance
- Depth, perspective, motion
- …

# Sub-challenge 5c: Model Induction

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

**Model Induction**



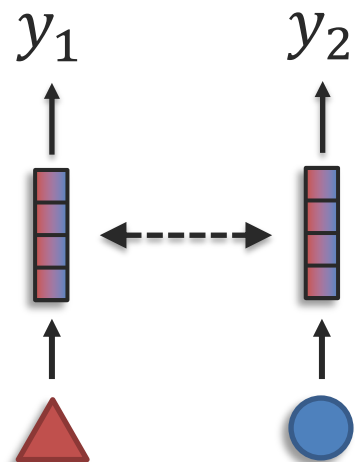$$y_1 \qquad y_2$$

# Sub-challenge 5c: Model Induction

Language △

Visual ●
(image)

**Information in both modalities = Y**
- Described people, objects, actions
- Illustrative gestures, motion

$y_1$    $y_2$

Common behavior: $X_1 \perp X_2 \mid Y$.
Or equivalently: $I(X_1; X_2 \mid Y) = 0$.

**Multi-view redundancy assumption**

# Co-training

**Setup**

Common behavior: $X_1 \perp X_2 \mid Y$.
Or equivalently: $I(X_1; X_2 \mid Y) = 0$.

$y_1$     $y_2$

**Multi-view redundancy assumption**
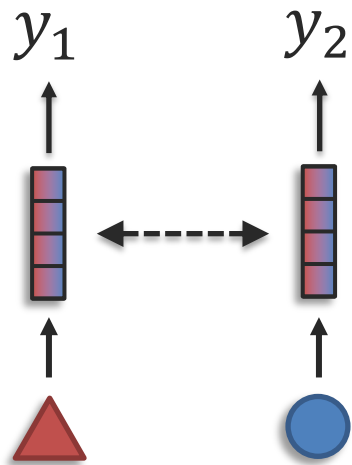
1. $X_1$ = text on the web page.
2. $X_2$ = text on hyperlinks pointing into the web page.
3. $Y$ = category of web page: academic, sports, news, music etc.

[Blum and Mitchell, Combining Labeled and Unlabeled Data with Co-Training. COLT 1998]

# Co-training

**Algorithm**

$y_1$       $y_2$

Assume:
1. Labeled data $\{X_1, X_2, Y\}$.
2. Unlabeled data $\{X_1, X_2\}$.

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

[Blum and Mitchell, Combining Labeled and Unlabeled Data with Co-Training. COLT 1998]

# Self-training

**Warmup: a single view – Self-training**

$y$

Assume:
1. Labeled data $\{X_1, Y\}$.
2. Unlabeled data $\{X_1\}$.

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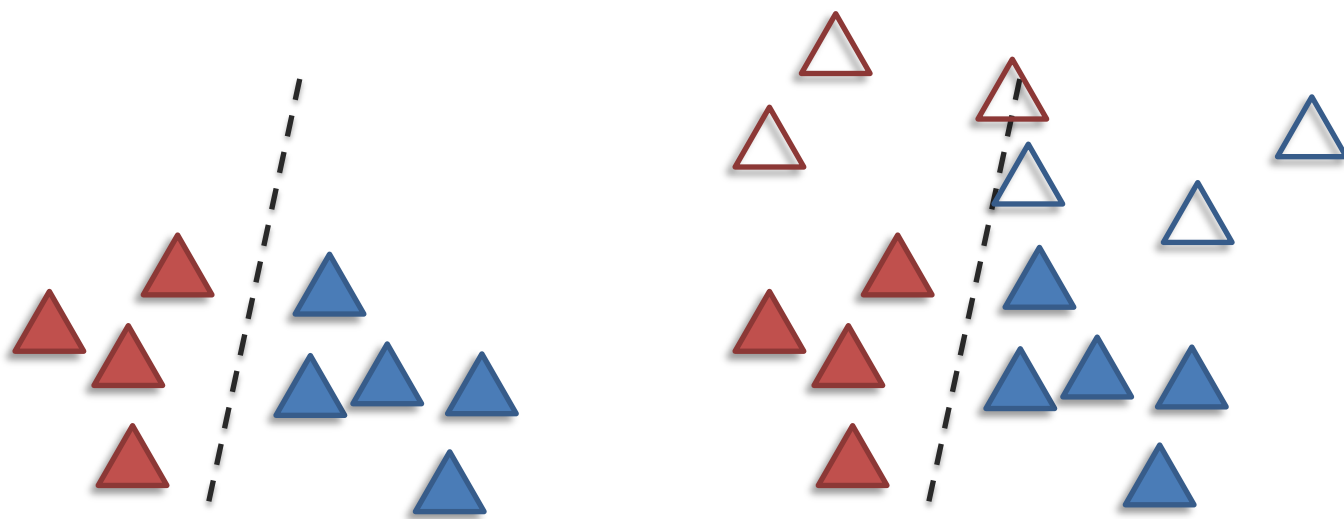
Test:
1. For a new unlabeled sample $\{X_1\}$, output $f_1(X_1)$.

# Self-training
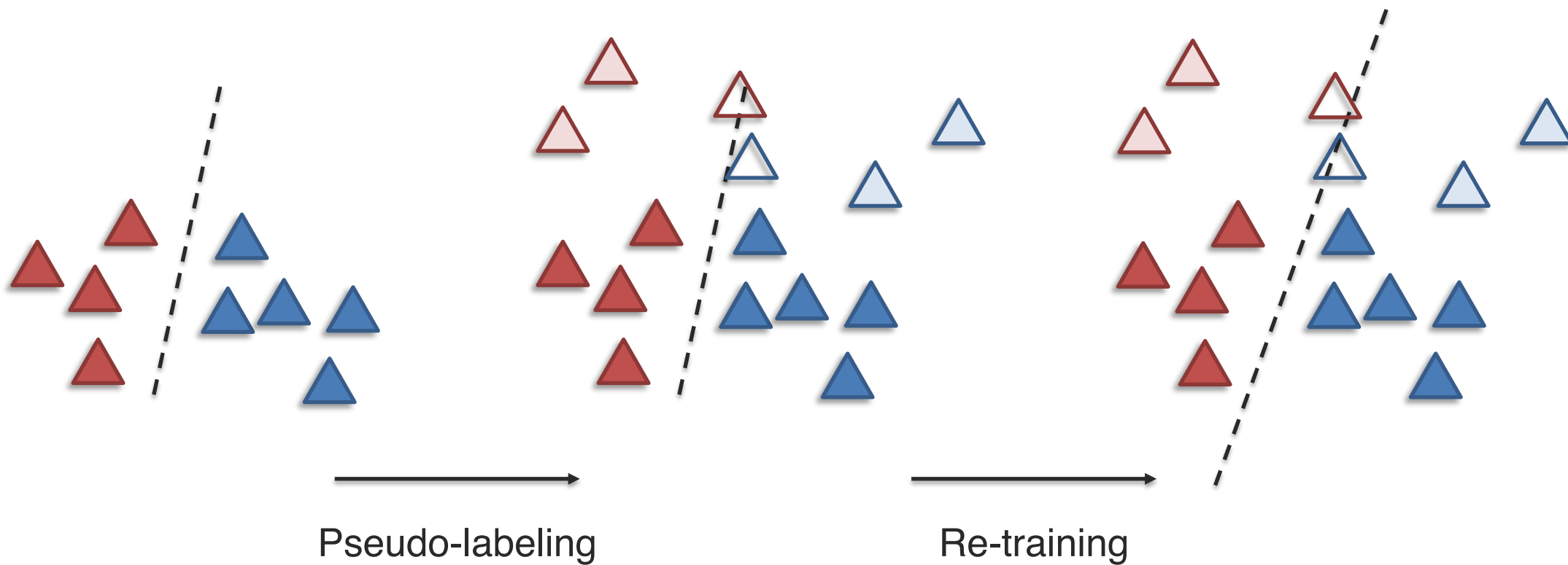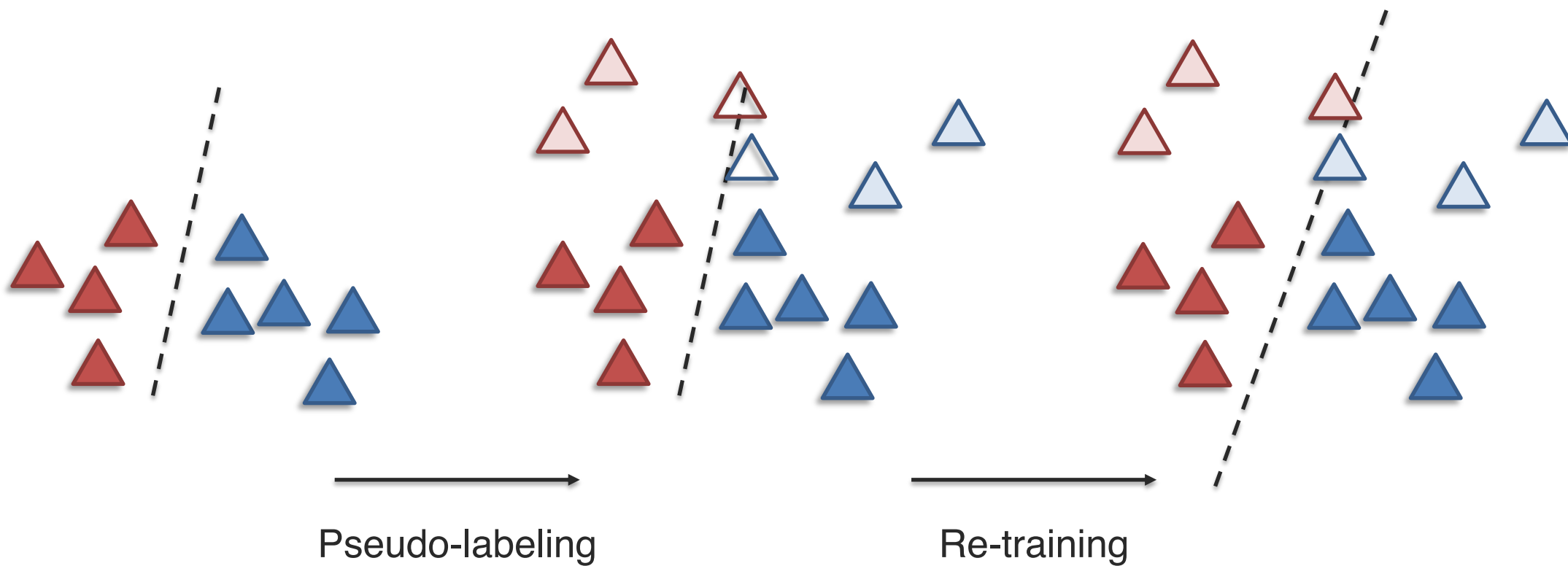
**Warmup: a single view – Self-training**

# Self-training

Warmup: a single view – Self-training



Pseudo-labeling

Re-training

# Self-training

**Warmup: a single view – Self-training**
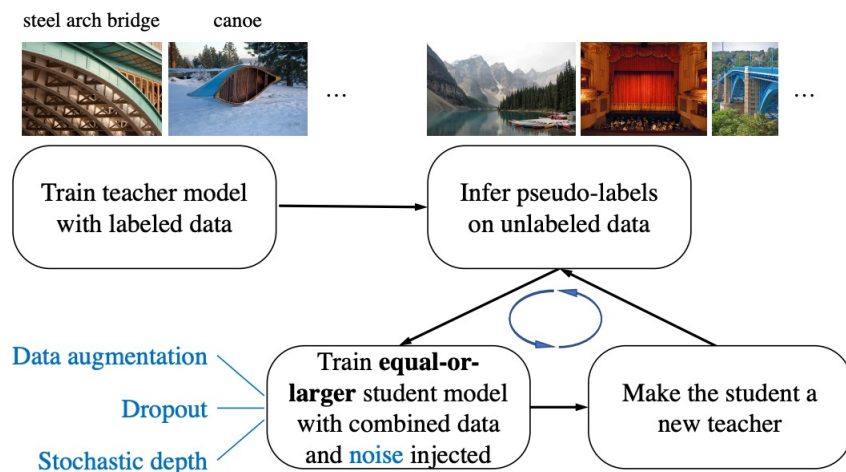


Pseudo-labeling

Re-training

# Self-training

**Key-words: semi-supervised learning, label propagation, domain adaptation/shift**

Critical:
1. Can't label all unlabeled data in one step, or you recover original classifier just trained on labeled data.
2. Sequence of pseudo-labeling is important to gradually shift classification boundary.
3. Input consistency regularization: shape of data space is important – implicit assumption that similar datapoints have similar labels (i.e., label consistency)



Input consistency:
- Data augmentation
- Adding noise

[Wei et al., Theoretical Analysis of Self-Training with Deep Networks on Unlabeled Data. ICLR 2021]

# Co-training

**Co-training**

$y_1$　　　$y_2$

Assume:
1. Labeled data $\{X_1, X_2, Y\}$.
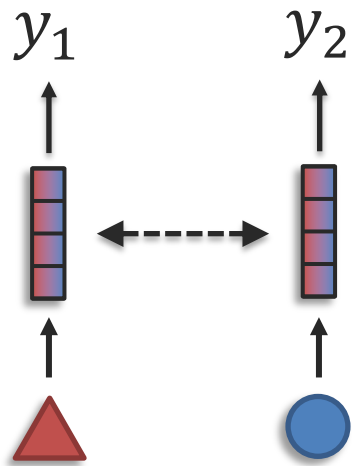2. Unlabeled data $\{X_1, X_2\}$.

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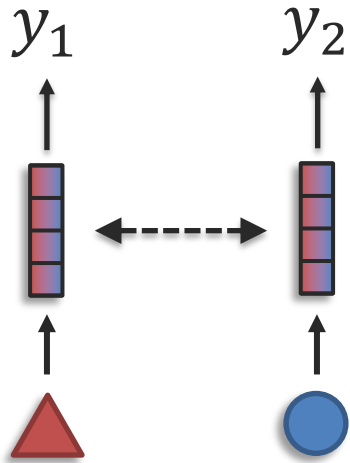4. Repeat until there are no more unlabeled samples.

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

[Blum and Mitchell, Combining Labeled and Unlabeled Data with Co-Training. COLT 1998]

# Co-training

**Co-training**

1. $X_1$ = text on the web page, $X_2$ = text on hyperlinks pointing into the web page.
3. $Y$ = category of web page: academic, sports, news, music etc.

$y_1$  $y_2$



**Louis-Philippe Morency**
Leonardo Associate Professor of Computer Science,
Language Technology Institute,
School of Computer Science, Carnegie Mellon University
Director, MultiComp Lab
Gates-Hillman Center (GHC) Office 5411,
5000 Forbes Avenue, Pittsburgh, PA 15213
Email: morency@cs.cmu.edu
Phone: (412) 268-5508

I am tenure-track Faculty at CMU Language Technology Institute where I lead the Multimodal Communication and Machine Learning Laboratory (MultiComp Lab). I was previously Research Faculty at USC Computer Science Department. I received my Ph.D. in Computer Science from MIT Computer Science and Artificial Intelligence Laboratory.

My research focuses on building the computational foundations to enable computers with the abilities to analyze, recognize and predict subtle human communicative behaviors during social interactions. Central to this research effort is the technical challenge of multimodal machine learning: mathematical foundation to study heterogeneous multimodal data and the contingency often found between modalities. This multi-disciplinary research topic overlaps the fields of multimodal interaction, social psychology, computer vision, machine learning and artificial intelligence, and has many applications in areas as diverse as medicine, robotics and education.

**Graduate Students Advising** (see all group members at MultiComp Lab website)

**Amir Ali Bagherzade**, Ph.D. program (LTI)
**Chaitanya Ahuja**, Ph.D. program (LTI)
**Volkan Cirik**, Ph.D. program (LTI co-supervised with Taylor Berg-Kirkpatrick)
**Alexandria Vail**, Ph.D. program (HCII)
**Paul Liang**, Ph.D. program (MLD, co-supervised with Ruslan Salakhutdinov)
**Hubert Tsai**, Ph.D. program (MLD, co-supervised with Ruslan Salakhutdinov)
**Torsten Wörtwein**, Ph.D. program (LTI)

Labeled, learn that '$X_1$ = CMU -> academic' and '$X_2$ = advised by -> academic'

**Paul Pu Liang**

Email: pliang(at)cs.cmu.edu
Office: Gates and Hillman Center 8011
5000 Forbes Avenue, Pittsburgh, PA 15213
Machine Learning Department and Language Technologies Institute, School of Computer Science, Carnegie Mellon University

[CV]  @pliang279  @pliang279  @lpwinniethepu

I am a fourth-year Ph.D. student in the Machine Learning Department at Carnegie Mellon University, advised by Louis-Philippe Morency and Ruslan Salakhutdinov. I also collaborate closely with Manuel Blum, Lenore Blum, and Daniel Rubin at Berkeley and Stanford. My research lies in the foundations of multimodal machine learning with applications in socially intelligent AI, understanding human and machine intelligence, natural language processing, healthcare, and education. As steps towards this goal, I work on:

Unlabeled, label using '$f_1: X_1$ = CMU -> academic' and learn that '$X_2$ = PhD program -> academic'

Another student -> Unlabeled, label using '$f_2: X_2$ = PhD program -> academic' and learn that '$X_1$ = Berkeley -> academic'

[Blum and Mitchell, Combining Labeled and Unlabeled Data with Co-Training. COLT 1998]

# Co-training

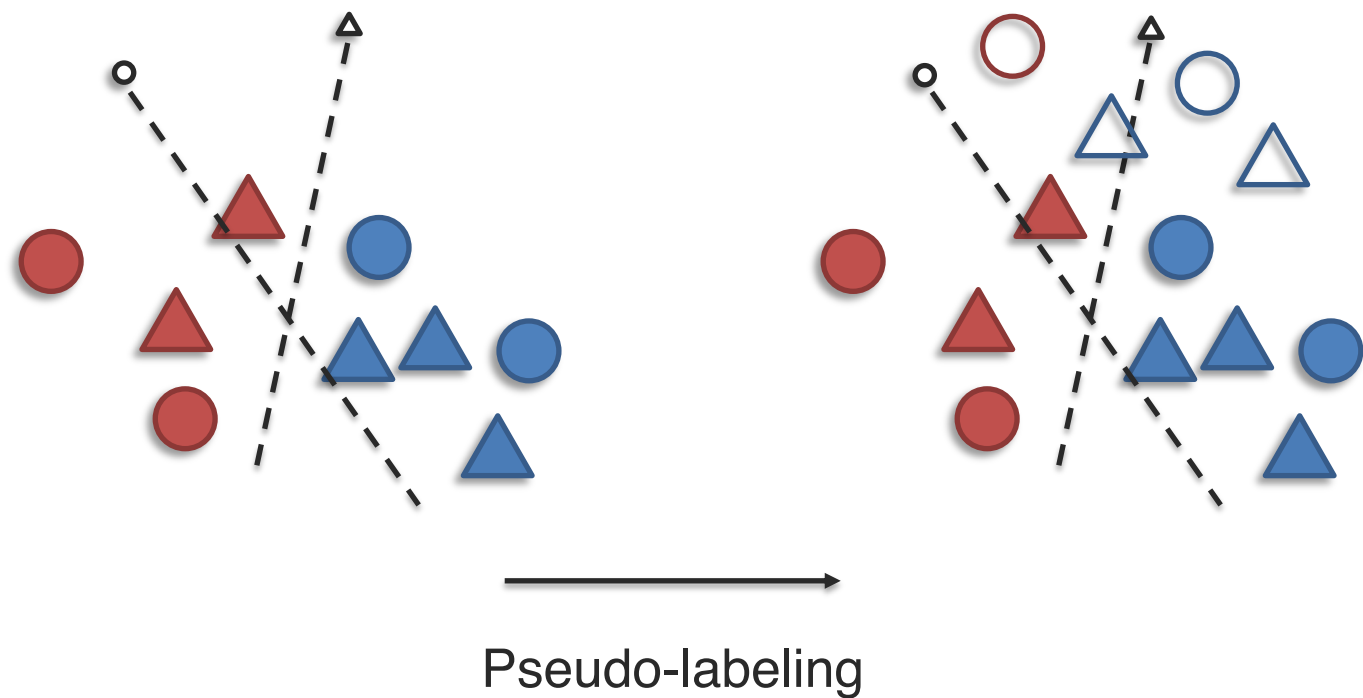**From self-training to co-training**

Assumptions:
1. Either view is sufficient to predict the label alone.
2. Views should be as independent as possible: examples where $f_1$ has high confidence but not $f_2$ and vice-versa.

[Blum and Mitchell, Combining Labeled and Unlabeled Data with Co-Training. COLT 1998]

# Co-training

**From self-training to co-training**



Pseudo-labeling

[Blum and Mitchell, Combining Labeled and Unlabeled Data with Co-Training. COLT 1998]

# Co-training

**From self-training to co-training**



Pseudo-labeling
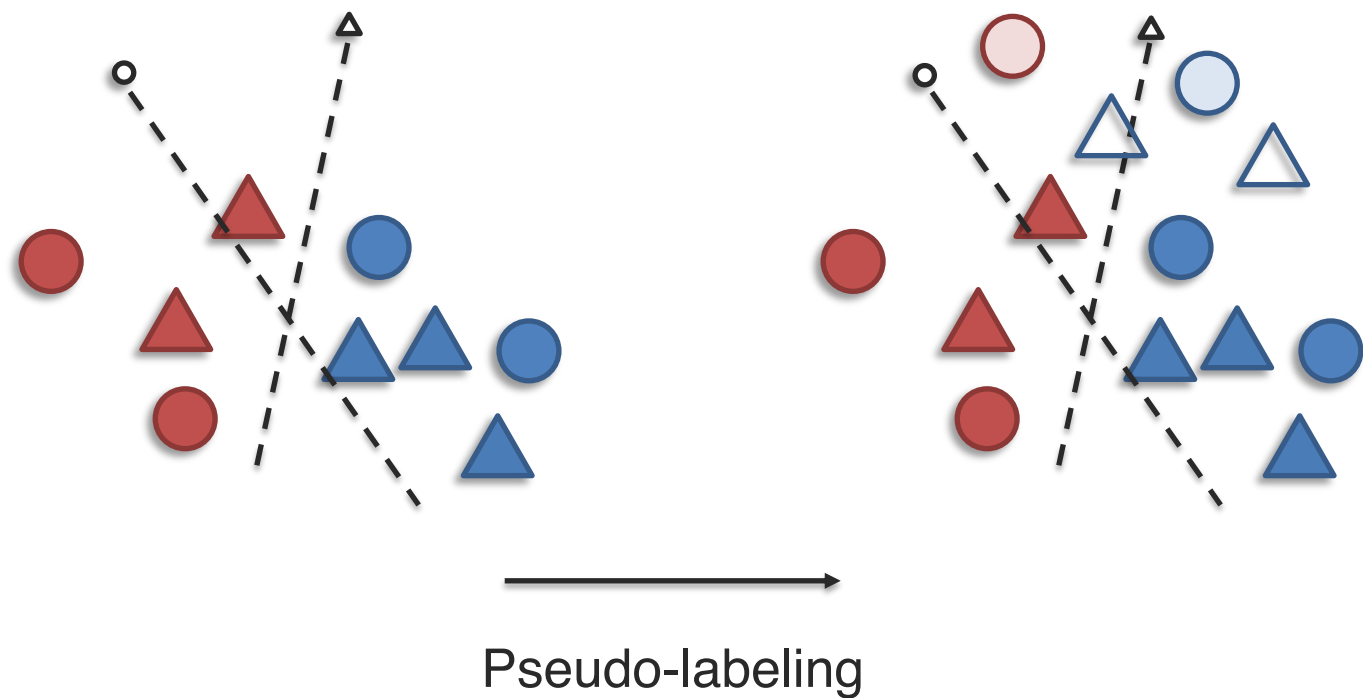
[Blum and Mitchell, Combining Labeled and Unlabeled Data with Co-Training. COLT 1998]

# Co-training

**From self-training to co-training**
**Key idea: functions on both views must be compatible and agree**



Pseudo-labeling

Re-training

[Blum and Mitchell, Combining Labeled and Unlabeled Data with Co-Training. COLT 1998]
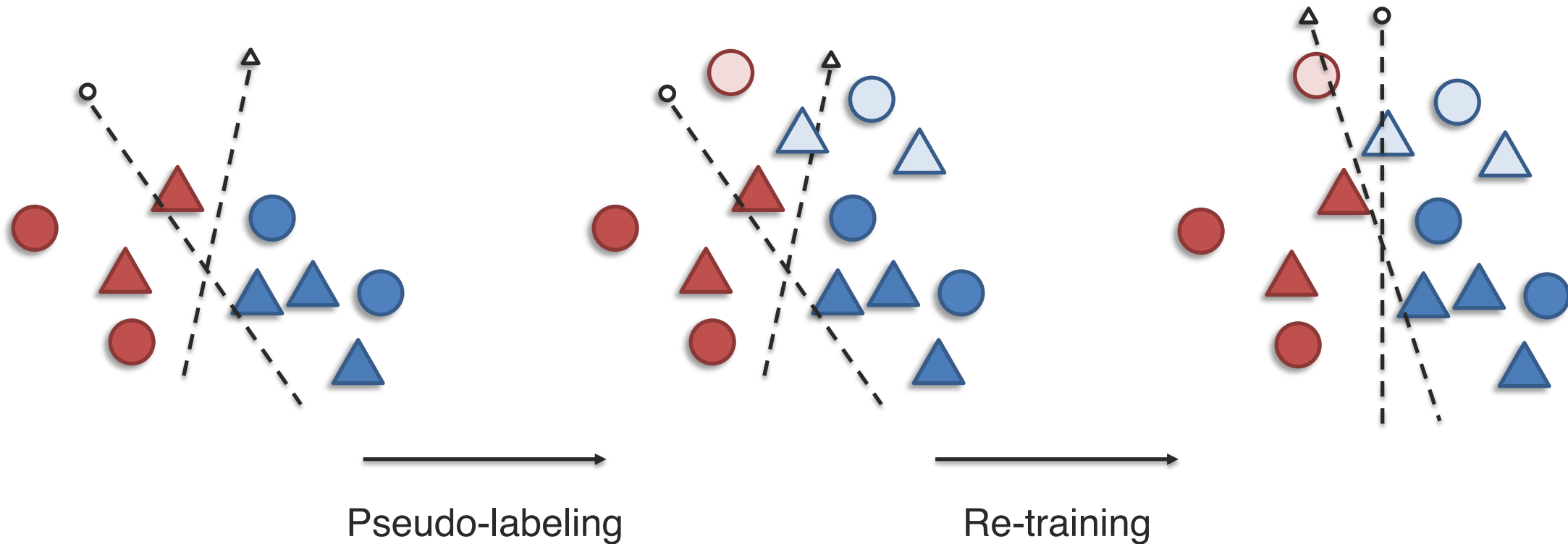
# Co-training

**From self-training to co-training**
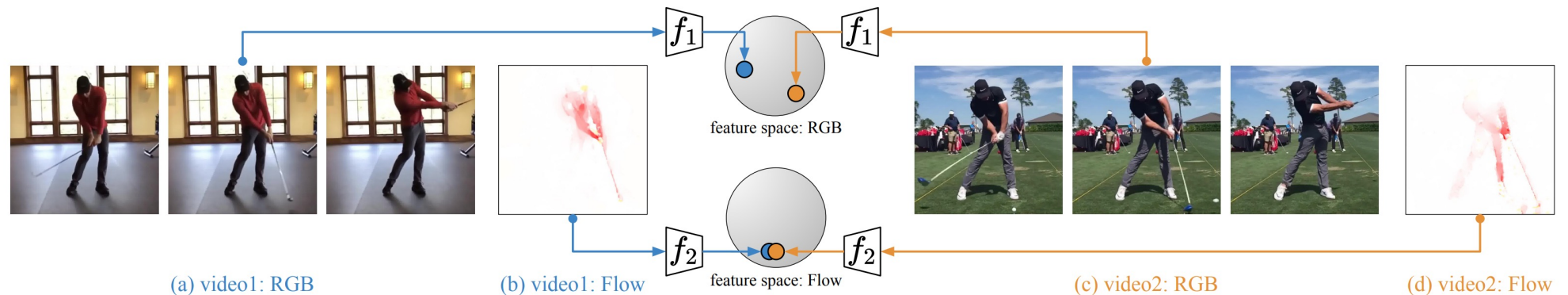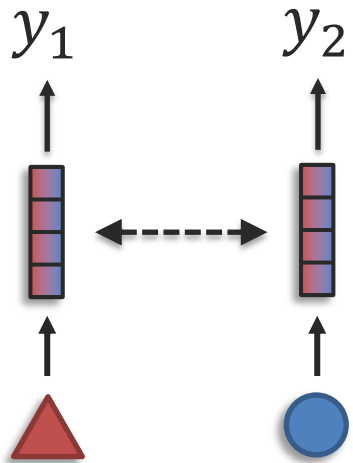**Key idea: functions on both views must be compatible and agree**

Intuitions:
1. Either view is sufficient to predict the label alone.
2. Views should be as independent as possible: examples where $f_1$ has high confidence but not $f_2$ and vice-versa.
3. Input consistency regularization: shape of data space is important – implicit assumption that similar datapoints have similar labels (i.e., label consistency).
→ In co-training, data from another view help us to supplement the label space!
→ Views independent given label = points in different views being in different spaces.
→ Both views must agree = input consistency which enables cross-view pseudo-labeling.
4. Eventually, will converge on 2 classifiers that agree and each separate both views.

[Blum and Mitchell, Combining Labeled and Unlabeled Data with Co-Training. COLT 1998]

# Co-training

**Recent applications of co-training**

$y_1$  $y_2$

Self-supervised learning with positive and negative samples
→ 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).
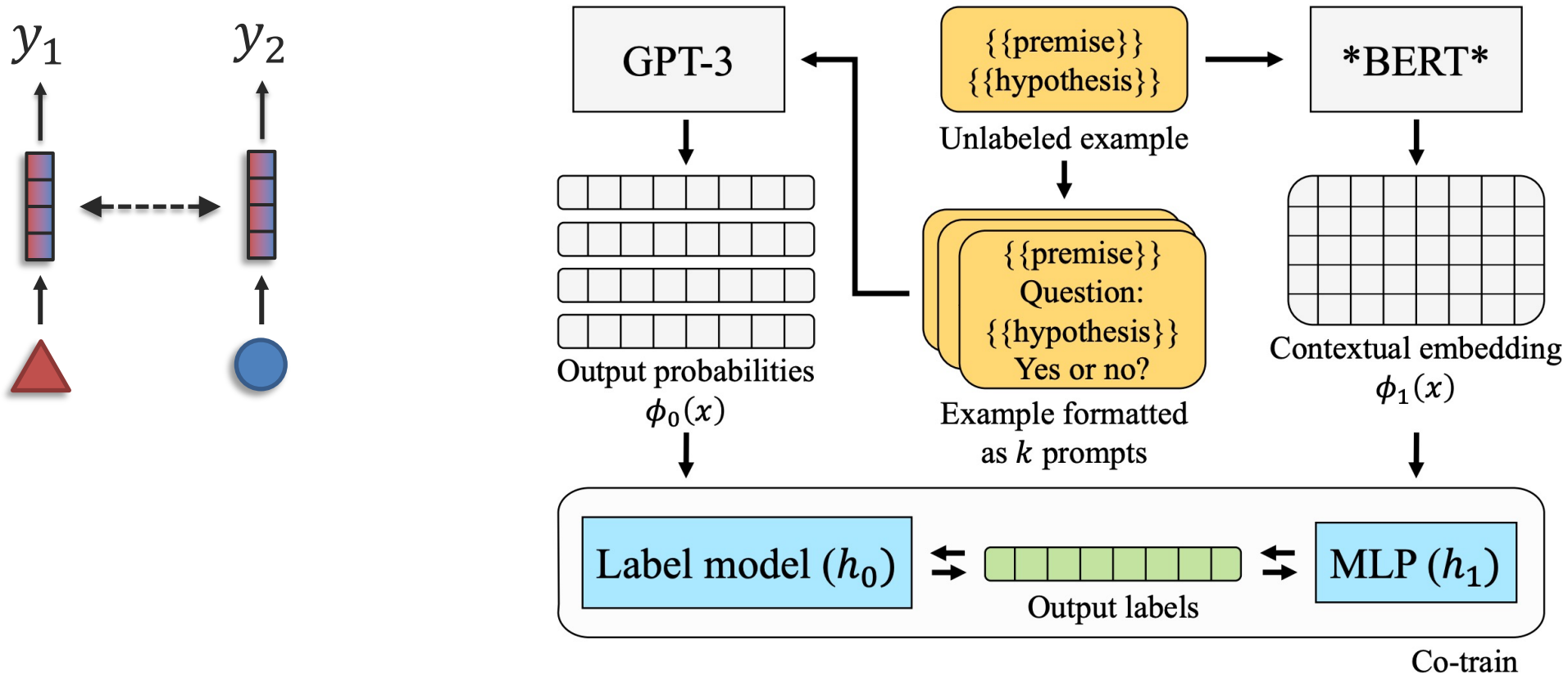→ Can use co-training between 2 RGB and flow contrastive learning modules.

$f_1$ feature space: RGB $f_1$

$f_2$ feature space: Flow $f_2$

(a) video1: RGB      (b) video1: Flow      (c) video2: RGB      (d) video2: Flow

[Han et al., Self-supervised Co-training for Video Representation Learning. NeurIPS 2020]

# Co-training

## Recent applications of co-training

Language-model prompting



[Lang et al., Co-training Improves Prompt-based Learning for Large Language Models. ICML 2022]
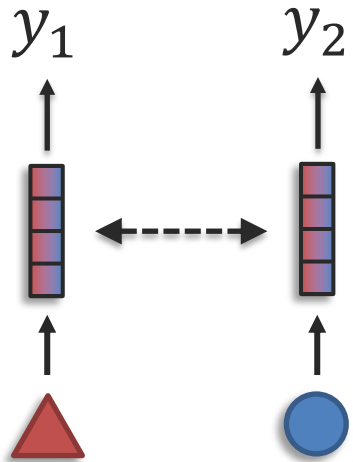
# Co-Regularization

**Co-regularization**

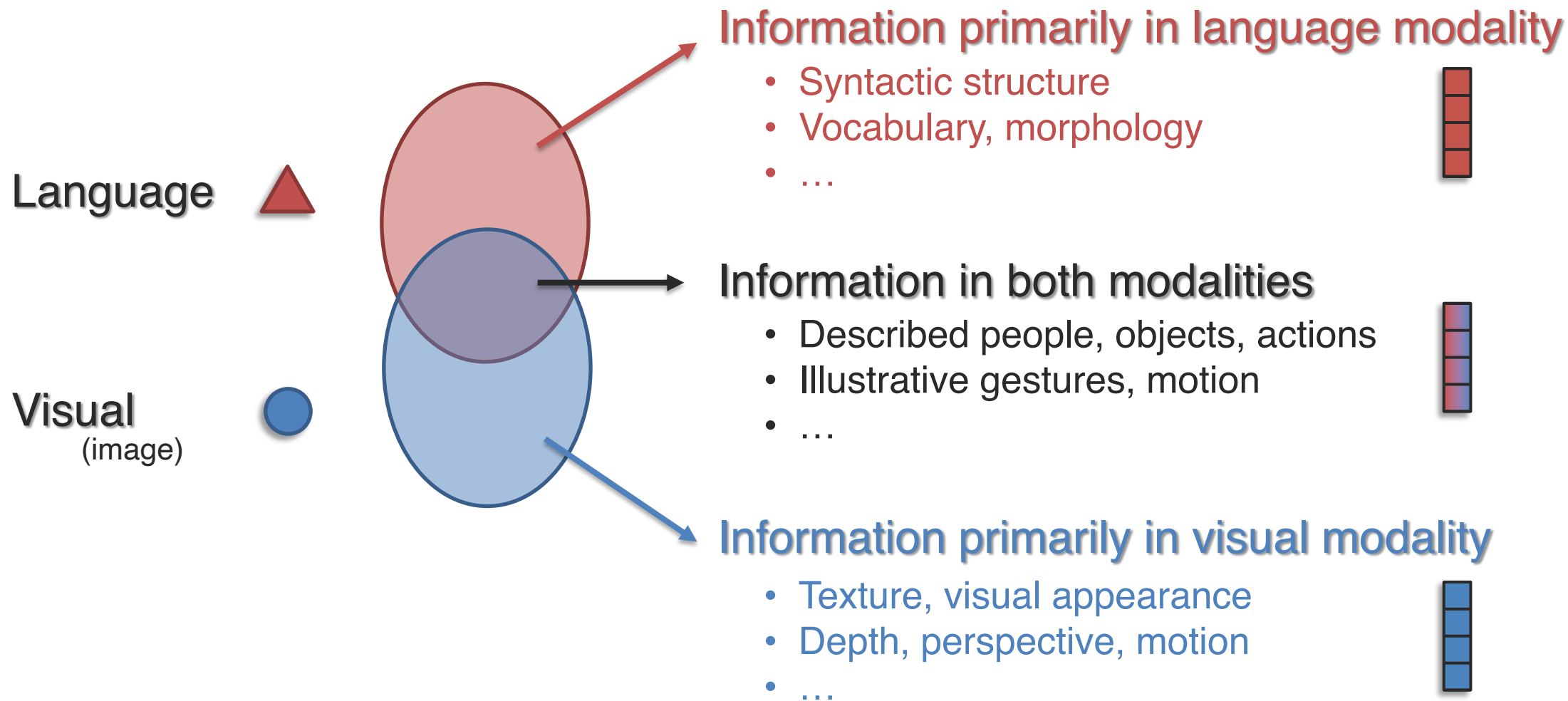Add a loss term to ensure both model predictions are similar:
$$L = (f_1(X_1) - f_2(X_2))^2$$

$y_1$  $y_2$

Recall representation coordination.

[Sridharan and Kakade, An Information Theoretic Framework for Multi-view Learning. COLT 2008]

# Sub-challenge 5c: Model Induction



**Information primarily in language modality**

- Syntactic structure
- Vocabulary, morphology
- …

**Information in both modalities**

- Described people, objects, actions
- Illustrative gestures, motion
- …

**Information primarily in visual modality**

- Texture, visual appearance
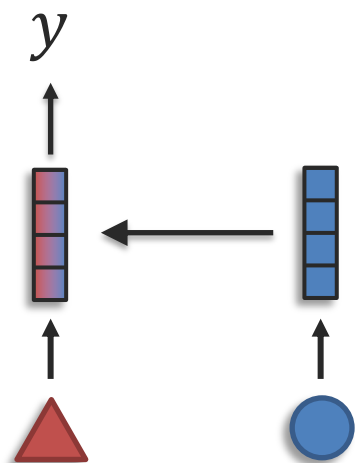- Depth, perspective, motion
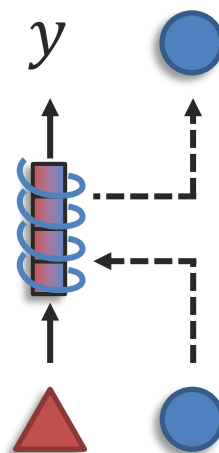- …

Language

Visual
(image)

# Summary: Transference

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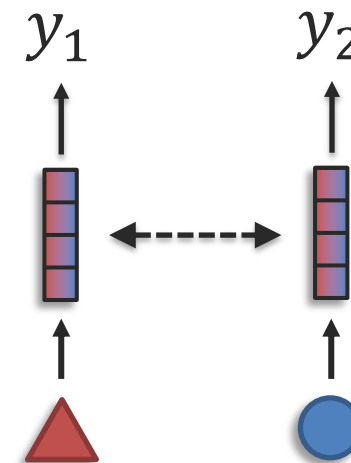
**Sub-challenges:**



**Transfer**      **Co-learning**      **Model Induction**

# More Transference

**Many more dimensions of transfer:**
→ **Multimodal {multitask, transfer, few-shot, meta} learning.**
→ **Domain adaptation, domain shift, label shift.**
→ **Core: representation, alignment, reasoning!**


**Open challenges:**
- Low-resource: little downstream data, lack of paired data, robustness (next section).
- Settings where SOTA unimodal encoders are not deep learning e.g., tabular data.
- Evaluating reasoning and robustness and large models.
- Limits of transfer beyond redundancy/joint information.
- Interpretability (next section).