

# Progressive Cross-Modal Causal Intervention for Long-Term Action Recognition

Shaowu Xu<sup>1</sup>, Xibin Jia<sup>1\*</sup>, Chao Fan<sup>1</sup>, Junyu Gao<sup>2</sup>, Jing Chang<sup>3</sup>, Qianmei Sun<sup>3</sup>  
 1 College of Computer Science, Beijing University of Technology, China; 2 Institute of Automation, Chinese Academy of Sciences, China; 3 Beijing Chao-yang Hospital, Capital Medical University, China



北京工业大学  
BEIJING UNIVERSITY OF TECHNOLOGY

CVPR  
JUNE 3-7, 2026

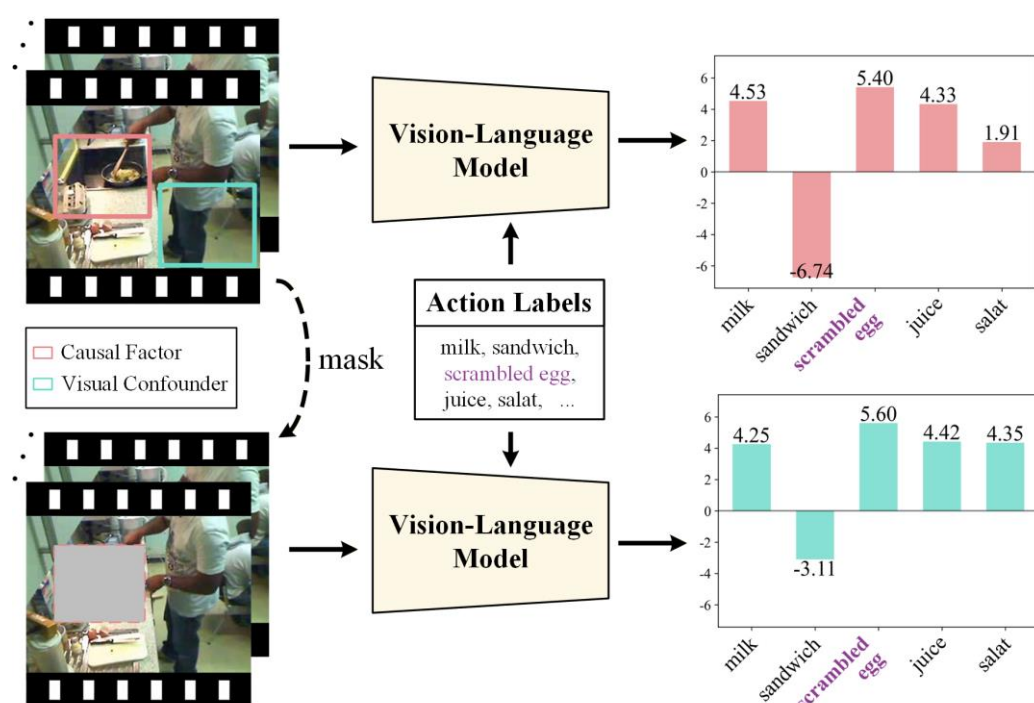


## Motivation: What Goes Wrong?

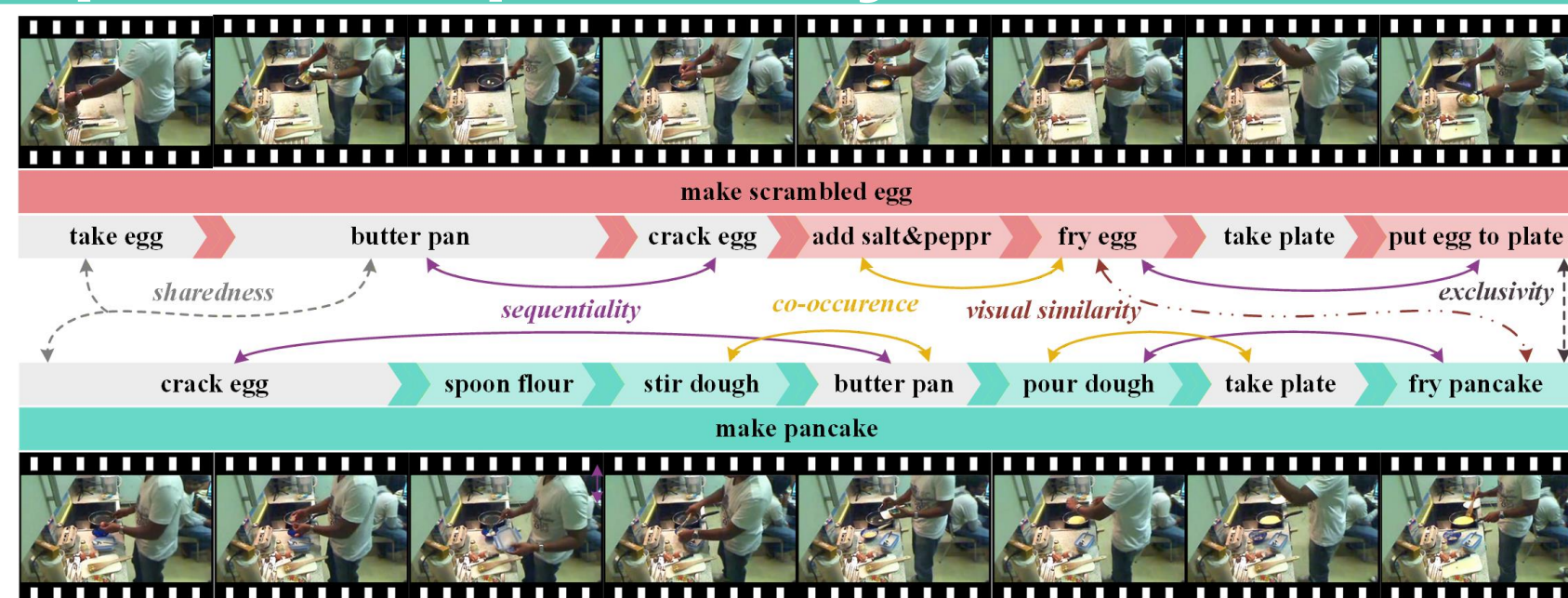
- LTAR must recognize long sequences of correlated atomic actions.
- VLM label supervision helps, but alignment can follow statistical shortcuts.
- **Three latent issues: co-occurrence hallucination (H), codependency illusion (I), visual confounders (C).**

## Example: Co-occurrence Hallucination

Masking causal regions can increase the score for "scrambled egg", showing reliance on non-causal visual cues.



## Example: Codependency Illusion



Isolated label modeling misses codependency cues among atomic actions, such as sharedness, exclusivity, etc.

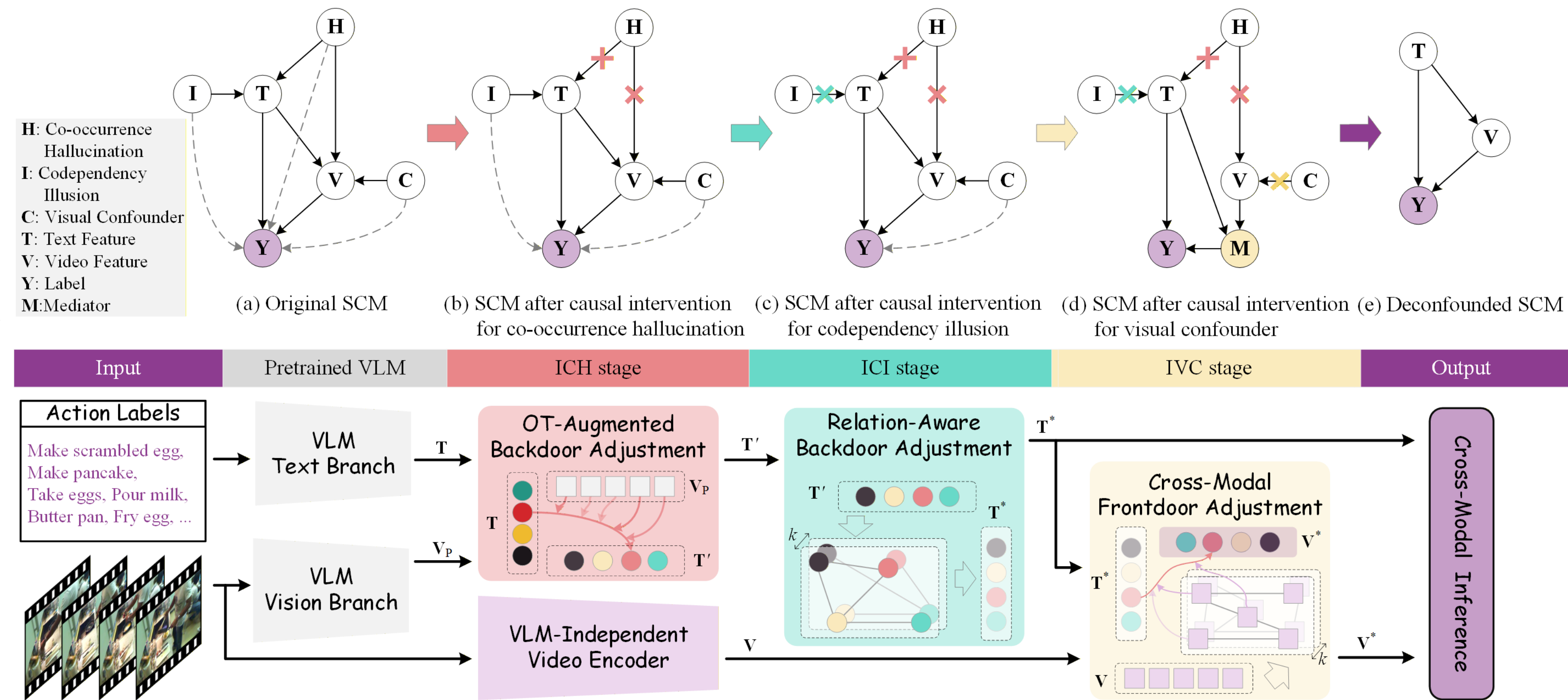
## Contributions

- A structural causal view of VLM-based LTAR with confounders H, I, and C.
- A progressive framework: ICH → ICI → IVC.
- SOTA performance on Breakfast, COIN, and Charades with lower inference cost.

**Core idea:** make text causal first, then use it as a mediator to deconfound vision.

## Method: Progressive Causal Intervention under SCM

PCMCi transforms a confounded VLM-based LTAR model into a deconfounded cross-modal model through Intervention for Co-occurrence Hallucination (ICH), Intervention for Codependency Illusion (ICI), and Intervention for Visual Confounder (IVC).



Latent confounders induce back-door paths:  $\{T, V\} \leftarrow H \rightarrow Y$     $T \leftarrow I \rightarrow Y$     $V \leftarrow C \rightarrow Y$    **SCM Paths**  
 Goal: predict Y from V and T while blocking spurious effects from H, I, and C.

### ICH Back-door adjustment for H

- VLM-independent encoder blocks  $H \rightarrow V$
- OT-guided back-door adjustment blocks  $H \rightarrow T$ , yielding  $T'$

$$S = \langle T, V^P \rangle$$

$$P^* = \operatorname{argmax}_{P \in \mathcal{U}} \langle P, -\log S \rangle + \lambda \mathcal{H}(P)$$

$$H = P^* \cdot V^P$$

$$P(Y | V, do(T)) \approx \sum_{h \in \mathcal{H}} P(Y | V, do(T), h) P(h) \approx P(Y | V, [T, H] W^H) = P(Y | V, T')$$

### ICI Back-door adjustment for I

- Relational cues estimates surrogate of I
- Back-door adjustment deconfounds  $T' \rightarrow T^*$ , blocking  $T \leftarrow I \rightarrow Y$

$$R(T') = \left( \bigoplus_{k=1}^K \mathcal{G}_{k(T'; \Theta_k)} \right) W^R$$

$$P(Y | V, do(T)) \approx \sum_{i \in \mathcal{I}} P(Y | V, do(T'), i) P(i) \approx P(Y | V, [T', R(T')] W^I) = P(Y | V, T^*)$$

### IVC Front-door adjustment for C

- Deconfounded  $T^*$  serves as mediator M
- Mediator-guided visual refinement deconfounds  $V \rightarrow V^*$ , blocking  $V \leftarrow C \rightarrow Y$

$$P(Y | do(V), T) = \sum_{v' \in \mathcal{V}'} P(Y | do(M), v', T) P(v')$$

$$= \sum_{m \in \mathcal{M}'} \sum_{v' \in \mathcal{V}'} P(Y | m, v', T) P(m | V) P(v')$$

$$\approx P(Y | [M, V'] W^C, T) = P(Y | V^*, T)$$

$$V' = \sigma(\mathcal{A}(M, \mathcal{F}(V); \Theta_M)) \mathcal{F}(V)$$

## Experiments

<b>97.46</b> Breakfast Acc	<b>90.51</b> Breakfast mAP	<b>86.54</b> COIN mAP	<b>53.3</b> Charades mAP
-------------------------------	-------------------------------	--------------------------	-----------------------------

PCMCi achieves competitive LTAR.

## Main Comparison

Method	B-Acc	B-mAP	C-Acc	C-mAP	FLOPs	Params
Text4Vis[1]	95.49	60.49	91.98	63.39	835	287
MA-LMM [2]	93.00	71.84	93.20	67.67	29063	7526
HierarQ [3]	97.18	76.32	<b>94.78</b>	70.10	37877	7881
<b>PCMCi</b>	<b>97.46</b>	<b>90.51</b>	94.53	<b>86.54</b>	<b>650</b>	<b>211</b>

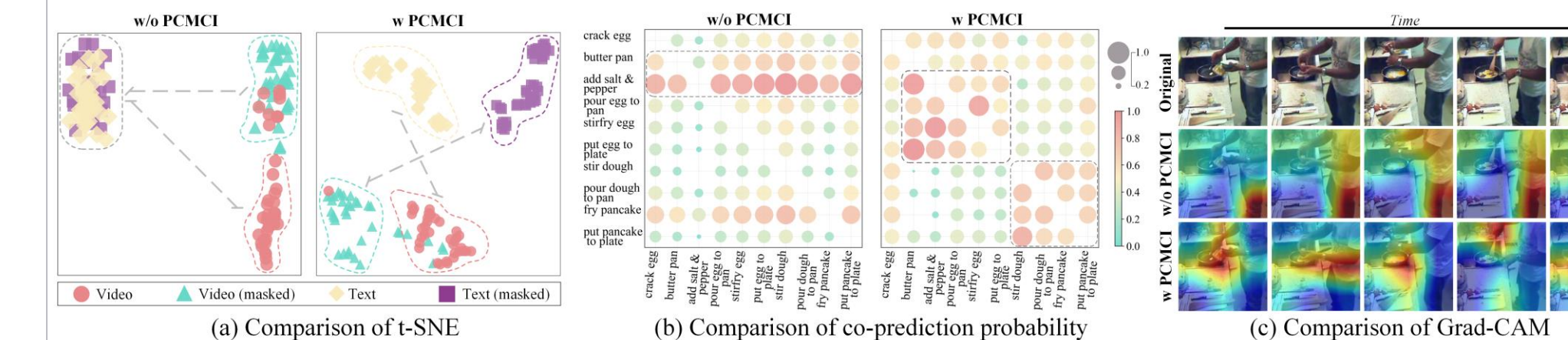
Breakfast (B) and COIN (C) results. PCMCi gains much higher mAP with far fewer FLOPs/parameters than LLM-based baselines.

## Ablation + Order

Setup	ICH	ICI	IVC	B-Acc	B-mAP
Base	-	-	-	91.55	80.23
ICH+ICI	✓	✓	-	95.77	85.59
IVC only	-	-	✓	93.24	88.31
<b>PCMCi</b>	✓	✓	✓	<b>97.46</b>	<b>90.51</b>

	Intervention Order	B-Acc	B-mAP
Order 1	IVC → ICH → ICI	92.96	81.76
Order 2	IVC → ICI → ICH	92.39	80.94
Order 3	ICI → IVC → ICH	94.08	83.13
Order 4	ICI → ICH → IVC	95.49	87.43
Order 5	ICH → IVC → ICI	96.34	89.72
<b>PCMCi</b>	<b>ICH → ICI → IVC</b>	<b>97.46</b>	<b>90.51</b>

## Visualization Evidence



(a) less hallucinated alignment; (b) better action co-dependency; (c) attention shifts from attire/background to causal regions.