

TIPS Over Tricks: Simple Prompts For Effective Zero-Shot Anomaly Detection



A. Salehi¹, E. Karami¹, Se. Noey³, Sa. Noey³, M. Yamada², R. Hosseini¹, **M. Sabokrou^{2*}**

¹ University of Tehran, ² Okinawa Institute of Science and Technology, ³ Amirkabir University of Technology

*mohammad.sabokrou@oist.jp

PROBLEM & MOTIVATION

- Anomaly detection:** industrial inspection and medical imaging.
- Traditional methods require **normal** data, which is often not available (Privacy) or is expensive
- Zero-shot anomaly detection (ZSAD)** solves this by using vision-language models (VLMs) like **CLIP**.

⚠️ CLIP LIMITATION:

- 🎯 Poor spatial alignment
- 🔍 Limited anomaly sensitivity
- 🧩 Requires complex add-ons (adapters, attention tricks, ...)

✓ TIPS Solution:

- ✓ Better backbone
- ✓ Spatial alignment
- ✓ Simple prompts

KEY CONTRIBUTIONS

- Rethink the Backbone**
Use TIPS instead of relying on CLIP fixes.
- Decoupled Prompting**
Separate prompts for detection and localization.
- Global-Local Fusion**
Combine image-level and pixel-level evidence.
- Simple Zero-Shot Pipeline**
Lightweight design with no complex add-ons.

QUANTITATIVE RESULTS

Zero-shot performance

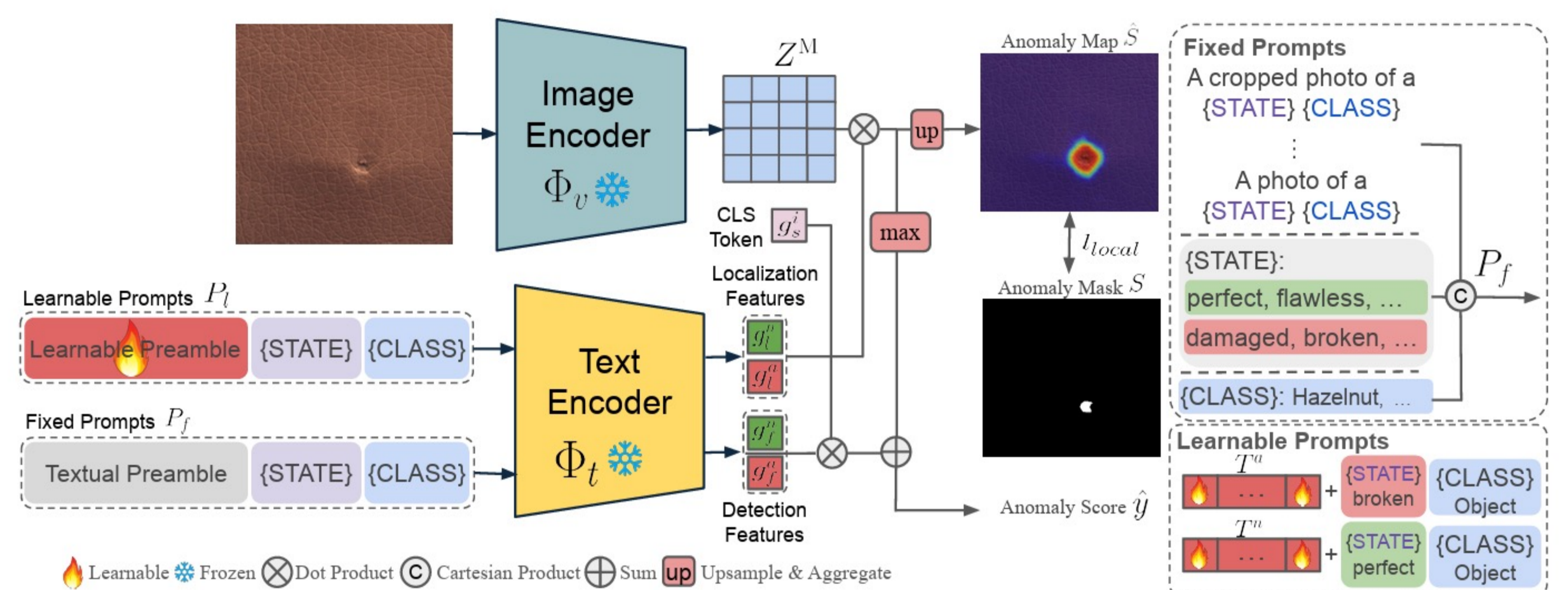
Setting	Metrics	Ours	Gain vs. Best Prior
Industrial Image-level	AUROC / AP / F1-max	92.5 / 93.1 / 90.7	+2.3 / +3.9 / +1.1
Industrial Pixel-level	AUROC / AUPRO / F1-max	96.5 / 90.2 / 51.7	+2.0 / +6.9 / +1.5
Medical Pixel-level	AUROC / AUPRO / F1-max	87.2 / 71.4 / 54.8	+3.2 / +4.4 / +5.3

Strong Performance Gains
Both anomaly detection and localization

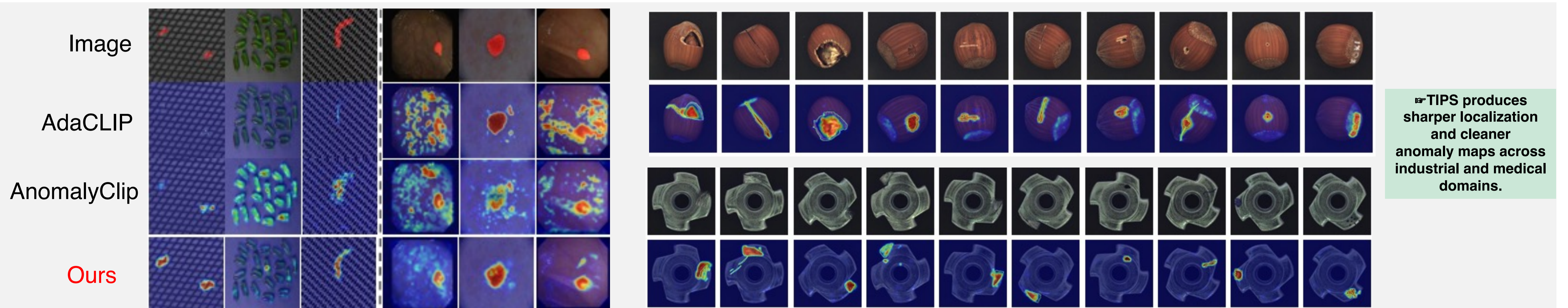
Average improvement over CLIP-based baseline

ARCHITECTURE

- Inputs:**
 - Image \rightarrow ViT \rightarrow global (g_s^i, g_o^i) + patch features (Z^M)
 - Prompts \rightarrow Transformer \rightarrow textual **prototypes**
- Scoring:**
 - Fixed prototypes (g_f) \rightarrow image-level scoring
 - Learnable prototypes (g_l) \rightarrow pixel-level anomaly map
 - Visual embeddings (e) $\rightarrow e \in \{g_s^i, g_o^i, Z^M\}$
$$p(e, G_{(\cdot)}) = \text{softmax}\left(\frac{G_{(\cdot)}^T e}{\tau}\right), \quad p_a = [p(e, G_{(\cdot)})]_a$$
- Combine:** global score + max(local anomaly)
- Output:**
 - anomaly score + anomaly map



QUALITATIVE RESULTS



CONCLUSION

- Using a better **backbone (TIPS)** is more effective than patching CLIP.
- Identified a **key issue**:
 - ⚠️ Global-local feature gap
- Solved with:
 - Decoupled prompting**
 - Local-global score fusion**

A simple, **spatially-aware** design can outperform complex CLIP-based adaptations in zero-shot anomaly detection.

PAPER & CODE

[arXiv Paper](#)



[GitHub Repo](#)

