

DISENTANGLING LATENT SHIFTS OF IN-CONTEXT LEARNING WITH WEAK SUPERVISION

Josip Jukić Jan Šnajder

TakeLab, University of Zagreb



NeurIPS 2025

- **In-Context Learning (ICL) is unstable:** performance varies with demonstration choice/order.
- **ICL is inefficient:** long prompts increase latency and cost; limited by context window.
- **Idea:** Treat ICL as weak supervision. Learn the *latent shift* induced by demos and **store it in an adapter**.

LINEARIZED ATTENTION VIEW

$$\mathbf{f}_{\text{AH}}(\mathbf{x}_q^{(t)}) \approx \underbrace{\mathbf{W}_{\text{ZS}} \mathbf{q}^{(t)}}_{\text{zero-shot}} + \underbrace{\Delta \mathbf{W}_{\text{ICL}} \mathbf{q}^{(t)}}_{\text{latent shift from demos}}$$

- Prior analyses often assume **linear attention**, neglecting nonlinear/residual dynamics.
- **Goal**: learn $\Delta \mathbf{W}_{\text{ICL}}$ into a compact adapter.

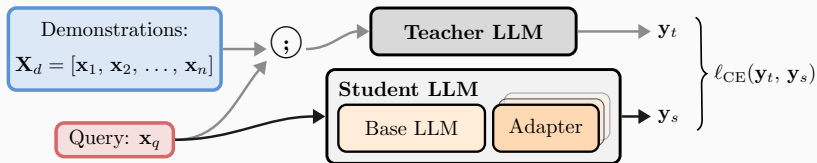
$\mathbf{x}_q^{(t)}$ – query token at step t

$\mathbf{q}^{(t)}$ – query vector

\mathbf{W}_{ZS} – zero-shot weights

$\Delta \mathbf{W}_{\text{ICL}}$ – demo-induced shift.

WEAK SUPERVISION



$$\mathcal{L} = \sum_{\mathbf{x}_q \in \mathcal{D}_{\text{unlab}}} \ell_{\text{CE}}(\mathbf{f}_{\text{teacher}}([\mathbf{X}_d; \mathbf{x}_q]), \mathbf{f}_{\text{student}}(\mathbf{x}_q))$$

- Teacher conditions on **demos + query**; student uses **query only**.
- Student adapter learns to **match teacher logits** \Rightarrow weakly supervised.

GENERALIZATION EXPERIMENTS

Method	GLUE						MMLU
	RTE	QNLI	MNLI	COLA	MRPC	QQP	MISC
<i>n</i> -shot	75.1 _{6.5}	77.0 _{5.5}	68.0 _{3.0}	58.5 _{4.0}	74.0 _{2.5}	70.0 _{3.0}	84.0 _{4.0}
PBFT	73.2 _{3.8}	77.8 _{6.0}	67.4 _{3.5}	56.5 _{3.0}	72.0 _{2.0}	68.0 _{2.5}	83.5 _{4.5}
Batch-ICL	77.8 _{4.7}	78.0 _{6.0}	70.9 _{3.5}	59.8 _{3.7}	75.2 _{2.2}	72.5 _{2.7}	81.0 _{2.5}
WILDA	86.0 _{0.6}	81.4 _{2.2}	73.1 _{2.0}	64.3 _{2.2}	77.7 _{1.5}	73.1 _{1.8}	88.0 _{2.2}

ID generalization accuracy for **Llama 3 (8B)** in the 16-shot setup.

- WILDA **achieves the strongest ID generalization**, outperforming standard ICL and related methods.
- **Highly stable**: variance is reduced compared to standard ICL.
- **Strong OOD generalization**: WILDA maintains high accuracy and low variance when evaluated on near-OOD GLUE pairs.

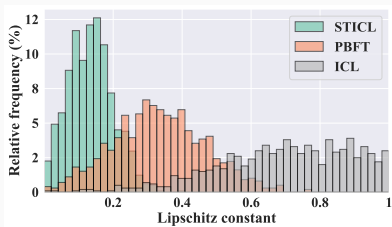
KNOWLEDGE FUSION VIA ADAPTER ARITHMETIC

Demos	Method	GLUE						MMLU
		RTE	QNLI	MNLI	COLA	MRPC	QQP	MISC
32	<i>n</i> -shot	75.3 _{3.2}	77.7 _{2.9}	69.1 _{1.9}	58.3 _{1.5}	76.4 _{2.2}	74.2 _{1.9}	84.5 _{2.1}
32	WILDA	87.9 _{0.6}	83.1 _{0.9}	74.0 _{1.1}	64.6 _{1.2}	79.4 _{0.6}	74.8 _{1.5}	89.0 _{0.4}
2 × 16	WILDA	87.1 _{1.6}	81.5 _{5.0}	75.5 _{2.5}	68.4 _{1.8}	78.5 _{1.4}	74.1 _{1.6}	89.5 _{2.0}

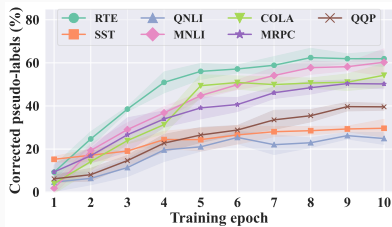
ID generalization accuracy for **Llama 3 (8B)** with fused demonstrations.

- **Adapter arithmetic** merges latent shifts from multiple subsets without retraining.
- **Stable fusion:** variance remains low as subsets increase.
- Enables **scalable task composition** – beyond context-window limits.

WEAK-TO-STRONG GENERALIZATION I



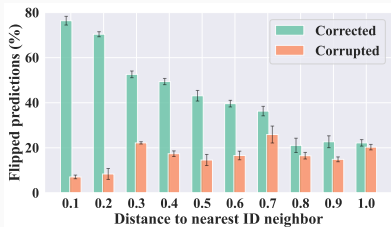
(a) Approximated Lipschitz constants



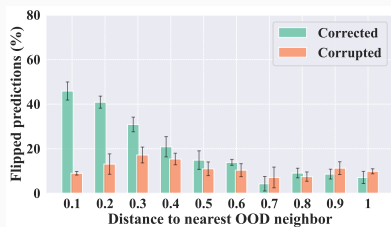
(b) Pseudo-label correction over epochs

- WILDA achieves strong local stability
- **Pseudo-label** corrections steadily increase across epochs

WEAK-TO-STRONG GENERALIZATION II



(c) ID corrected/corrupted rates



(d) OOD corrected/corrupted rates

- Correction rates fall sharply with distance to the nearest correctly pseudo-labeled neighbor → **coverage expansion**.
- The same pattern appears on OOD data → WILDA generalizes corrections beyond the ID manifold.