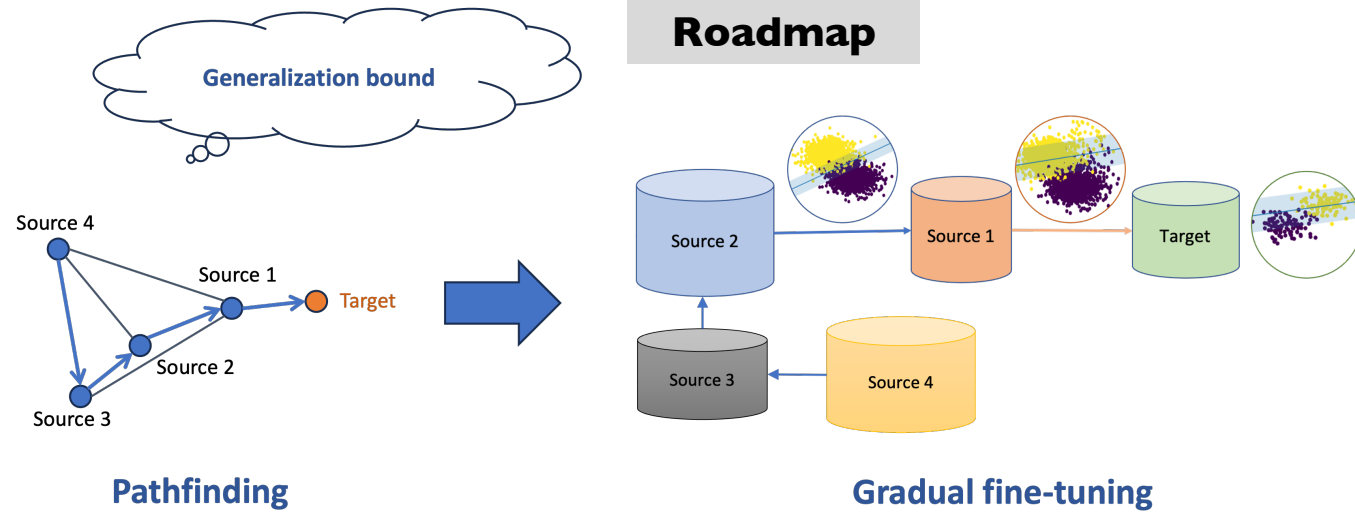


Takeaways

- Gradual fine-tuning (GFT) of models across multiple source domains to one target, represented as an undirected weighted graph.
- New generalization bound for GFT along any path in graph, guiding optimal training order.
- Lightweight graph-routing pathfinding strategies that minimize error bound and attain practical target accuracy efficiently.

Roadmap

In a nutshell:

- Derive GFT gen error bound.
- Find path to make small bound.
- Follow path to fine-tune model.

Problem definition

MSUDA adapts a model to a target domain using multiple source domains without data labels in the target domain.

Challenges: Lack of target domain access; Costly selection; Distant sources.

Objective: Minimize expected target risk by training on multiple sources

$$h_T^* = \arg \min_{h \in \mathcal{H}} \mathbb{E}_{\{x,y\} \in D_T} [\mathcal{L}(h(x), y)].$$

where we explore \mathcal{H} by source data.

GFT

At timestep t , classifier is trained on source S_t with size (magnitude) n_t , from initialization by the previous classifier \hat{h}_{t-1}

$$\hat{h}_t \leftarrow \arg \min_{h \in \mathcal{H}, h^0 \leftarrow \hat{h}_{t-1}} \frac{1}{n_t} \sum_{\{x,y\} \in S_t} \mathcal{L}(h(x), y).$$

Results on MultiNLI and Sentiment Analysis

Table 1: Accuracy comparison on 5 target domains from the MultiNLI dataset, in mean \pm std. Subscripts of average accuracies denote relative decreases to the best performance. Repeated experiments are conducted above identical set of seeds for training.

Method	Target Domain					Avg Acc.
	Fiction	Government	Telephone	Slate	Travel	
ALL SOURCES	76.62 \pm 0.67	72.34 \pm 1.57	71.94 \pm 1.37	71.09 \pm 1.12	72.47 \pm 2.02	72.89 (\downarrow 4.7%)
CLOSEST	74.97 \pm 0.34	72.88 \pm 1.19	72.24 \pm 0.71	73.50 \pm 1.28	71.36 \pm 0.85	72.99 (\downarrow 4.6%)
SEAL-SHAP Xu et al. (2021)	74.70 \pm 1.62	75.39 \pm 0.75	74.63 \pm 2.05	73.37 \pm 0.69	75.70 \pm 3.07	74.75 (\downarrow 2.3%)
TGFT	77.43 \pm 1.78	77.19 \pm 2.13	72.89 \pm 2.08	74.35 \pm 1.59	74.68 \pm 4.43	75.30 (\downarrow 1.6%)
NNGFT	78.03 \pm 2.34	76.95 \pm 2.14	73.74 \pm 2.19	77.03 \pm 6.27	76.76 \pm 2.19	76.50 (0.0%)
SPGFT	76.40 \pm 1.31	73.91 \pm 5.31	73.05 \pm 1.69	71.00 \pm 3.39	73.14 \pm 2.85	73.50 (\downarrow 3.9%)
MSTGFT	76.18 \pm 4.39	73.91 \pm 5.31	73.05 \pm 1.69	71.00 \pm 3.39	73.14 \pm 2.85	73.45 (\downarrow 4.0%)

Table 3: Accuracy comparison on 4 distant domains from the multi-domain sentiment analysis dataset, in mean \pm std. Subscripts of average accuracies denote relative decreases to the best performance. Repeated experiments are conducted above identical set of seeds for training.

Efficacy in distant domains

Method	Target Domain				Avg Acc.
	Books	Music	Electronics	Grocery	
ALL SOURCES	89.71 \pm 0.31	88.83 \pm 0.67	87.75 \pm 0.56	88.66 \pm 0.53	88.73 (\downarrow 0.5%)
CLOSEST	89.69 \pm 0.41	88.98 \pm 0.22	83.66 \pm 0.78	88.62 \pm 0.88	87.73 (\downarrow 1.6%)
SEAL-SHAP Xu et al. (2021)	84.85 \pm 1.80	85.91 \pm 1.52	88.18 \pm 0.47	84.21 \pm 1.50	85.79 (\downarrow 3.8%)
NNGFT	89.33 \pm 0.53	89.85 \pm 0.32	87.65 \pm 0.04	89.85 \pm 0.82	89.17 (0.0%)

- 2.3%** and **3.9%** relative accuracy improvements over domain scoring SOTA, SEAL-Shap.
- Outperforms in-/out-domain gradual shift [Xu et al., 2021], which trains models with target data.

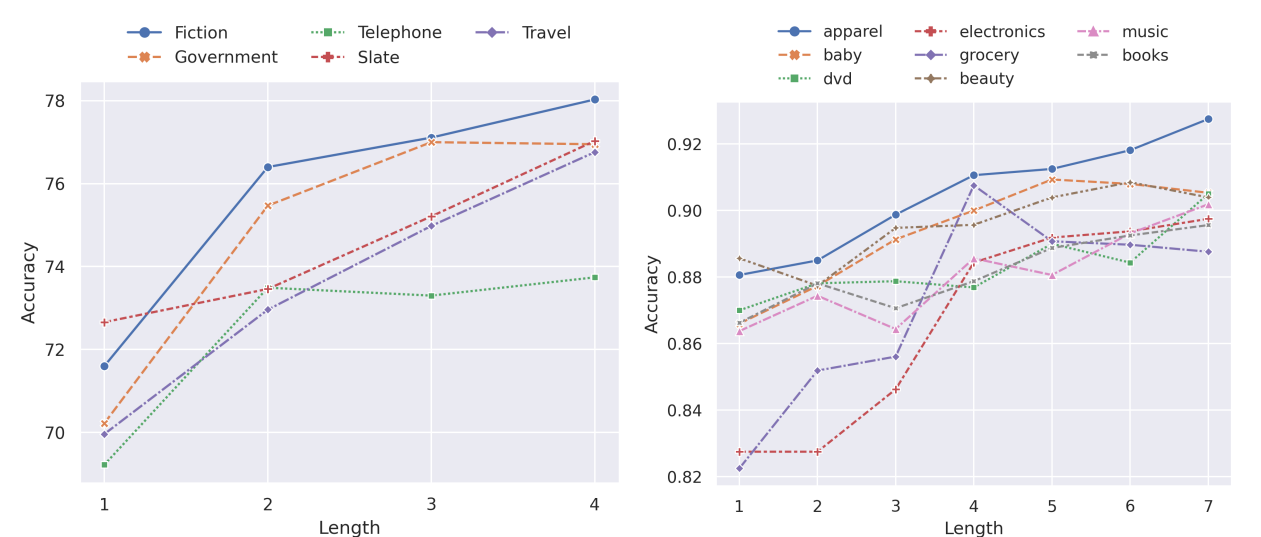
Additional setup info

Dataset: Amazon Review for SA.

Model: BERT-adapted.

Baselines: Trained on

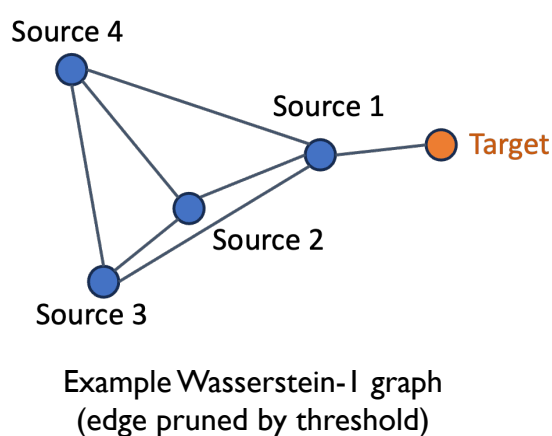
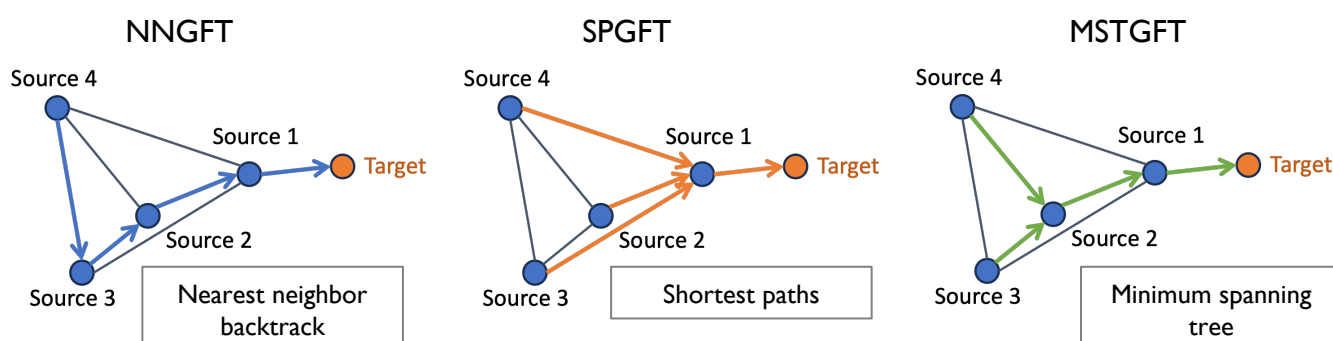
- all sources combined;
- closest source only;

Ablation study

- Indicates error bound $\Delta(\rho)$ takes more advantage from path magnitude $\text{mag}(\rho)$.
- Justifies performance tradeoffs for SPGFT and MSTGFT, which sacrifice path magnitude for length efficiency and thereby training speed.

Graph routing

- Low-cost source domain selection.
- Efficient use of distant domains.

**Graph routing paradigm**

- For each source domain S_i , find the most length-efficient path to target T .
- Select the one π^* with maximum magnitude in optimal paths P^* .

$$\pi^* \approx \arg \max_{\pi \in P^*} \text{mag}(\pi),$$

$$\text{where } P^* = \left\{ \arg \min_{\rho \in P_G(S_i, T)} \Delta(\rho), i \in [K] \right\}.$$

Future work

- Gap between UDA generalization error and graph theory.
- Source combining + pathfinding instead of source-by-source.
- Scalability to larger domain graphs.
- Exact / Stronger path optimality (not just error bounded).
- More expressive distance measure.



Join the
climb: Scan
to read!

