

A Task-Specific Transfer Learning Approach to Enhancing Small Molecule Retention Time Prediction with Limited Data

Yuhui Hong, Haixu Tang Luddy School of Informatics, Computing, and Engineering Indiana University, Bloomington, IN 47408, USA

The authors declare no competing financial interest.

Background | LC-MS retention time prediction



Liquid Chromatography-Mass Spectrometer



ψ

Background | LC-MS retention time prediction

METLIN-SMRT 80,038 small molecules ^[1]



Instrumentation	
LC System	Agilent 1100/1200 series
Mass Spectrometer	Quadrupole-time of flight (Q-TOF) G6538A (Agilent Technologies)
Column	Zorbax Extend-C18 reverse-phase (2.1 × 50 mm, 1.8 µm, Agilent Technologies)

Chromatographic Conditions	
Flow rate	100 μL/min
Mobile phase A	Water + 0.1% formic acid
Mobile phase B	Acetonitrile + 0.1% formic acid
Dead volume	40 μL
Dwell volume	900 µL



[1] Domingo-Almenara X, et al. Nat. Commun. 10, 5811 (2019)[2] Xue J, et al. Bioinformatics 40.3, btae084 (2024)



Compounds number (isomeric records) distribution of RT datasets in the RepoRT database ^[1] before preprocessing

Only 10 datasets \geq 300 compounds



[1] Domingo-Almenara X, et al. Nat. Commun. 10, 5811 (2019)[3] Kretschmer F, et al. Nat. Methods 21.2, 153-155 (2024)

.

Potential solution 1: Converting retention times to retention indices ^[4]







Potential solution 1: Converting retention times to retention indices ^[4]







Potential solution 1: X Converting retention times to retention indices







Potential solution 2: Transfer learning^[5, 6]





Upstream dataset, e.g. METLIN-SMRT^[5], LC-MSMS datasets^[6], etc.

[•]





Downstream dataset



[5] Kwon Y, et al. Anal. Chem. 95.47, 17273-17283 (2023)[6] Hong Y, et al. Bioinformatics 39.6, btad354 (2023)

Potential solution 2: Transfer learning ^[5]

Performance of Graph Isomorphism Network (GIN) on 45 datasets from the RepoRT database using METLIN-SMRT pre-training ^[3]





Potential solution 2: Transfer learning

X

- Large domain gap between datasets
- Limited downstream data



Upstream dataset, e.g. METLIN-SMRT



fine-tuning



Downstream dataset



Our solution:

Task-Specific Transfer Learning (TSTL)

Step 1: Joint datasets pre-training

Step 2: Integration of multiple fine-tuned models

Step 1: Joint datasets pre-training





Upstream dataset + Downstream dataset





Downstream dataset



Optimizer methods

 $\mathcal{L}(\emptyset)$

pre-trained model

best initialization for specific downstream task



Optimizer methods

$$\mathcal{L}(\emptyset) = l_t(\theta_t)$$

pre-trained model

fine-tuned model

best initialization for specific downstream task



Optimizer methods

$$\mathcal{L}(\emptyset) = l_t(\theta_t) = l_t(\emptyset - \alpha \nabla_{\emptyset} l_t(\emptyset))$$

inner loop: update θ_t initialized as \emptyset on downstream task

learning rate: α , β



Optimizer methods

learning rate: α , β

on downstream task



Step 2: Integration of multiple fine-tuned models





Results | data preprocessing

TL-difficult datasets (TL $R_2 < 0.8$)



Compounds number distribution of RepoRT database ^[3] after preprocessing



Results | diverse and correlation of up stream tasks



ψ

Results | performance of TSTL on retention time prediction





Results | performance of TSTL on retention time prediction



 $SC \Rightarrow TL-SMRT \rightarrow TL-ALL \rightarrow TSTL-SMRT \Rightarrow TSTL-ALL$



Takeaways

- We designed Task-Specific Transfer Learning (TSTL) for training neural networks with limited data
- TSTL incorporates 2 steps: joint pre-training and a greedy integration strategy
- Experiments demonstrate that TSTL outperforms TL on all TL-difficult datasets in RepoRT database

Future work

- Apply TSTL methodology to more predictions
- Enhance integration considering experimental condition correlations



Acknowledgement

We acknowledge the Center for Bioanalytical Metrology (CBM), an NSF Industry-University Cooperative Research Center, for providing funding under grant NSF IIP-1916645. This work was also partially supported by National Science Foundation grant DBI-2011271.

Visit other works from our lab!

TP 057: Al-driven visual proteomics for blood-based Alzheimer's disease biomarker discovery by LC-MS/MS and deep neural networks. *Qingyang Xiao et al.*

WP 451: Metabolite identification by spectral searching against predicted spectral library. *Chhavi Thakur et al.*





Manuscript is under preparing. Follow this GitHub repository for any updates!



References

[1] Domingo-Almenara, Xavier, et al. "The METLIN small molecule dataset for machine learning-based retention time matters: Transferable prediction of small molecule liquid prediction." Nature communications 10.1 (2019): 5811. [2] Xue, Jun, et al. "RT-Transformer: retention time prediction for metabolite annotation to assist in metabolite identification." Bioinformatics 40.3 (2024): btae084. [3] Kretschmer, Fleming, et al. "RepoRT: a comprehensive repository for small molecule retention times." Nature Methods 21.2 (2024): 153-155.

[4] Kretschmer, Fleming, et al. "Times are changing but order chromatography retention times." (2024). [5] Kwon, Youngchun, et al. "Retention time prediction through learning from a small training data set with a pretrained graph neural network." Analytical Chemistry 95.47 (2023): 17273-17283. [6] Hong, Yuhui, et al. "3DMolMS: prediction of tandem mass spectra from 3D molecular conformations." Bioinformatics 39.6 (2023): btad354.

