

Automatic Auxiliary Task Selection and Adaptive Weighting Boost Molecular Property Prediction

Zhiqiang Zhong^[1,2,3], Davide Mottin^[1]

[1]Department of Computer Science, Aarhus University

[2]Institute for Advanced Studies, University of Luxembourg

[3]Faculty of Science, Technology and Medicine, University of Luxembourg



INSTITUTE FOR ADVANCED STUDIES (IAS)



Co-funded by
the European Union

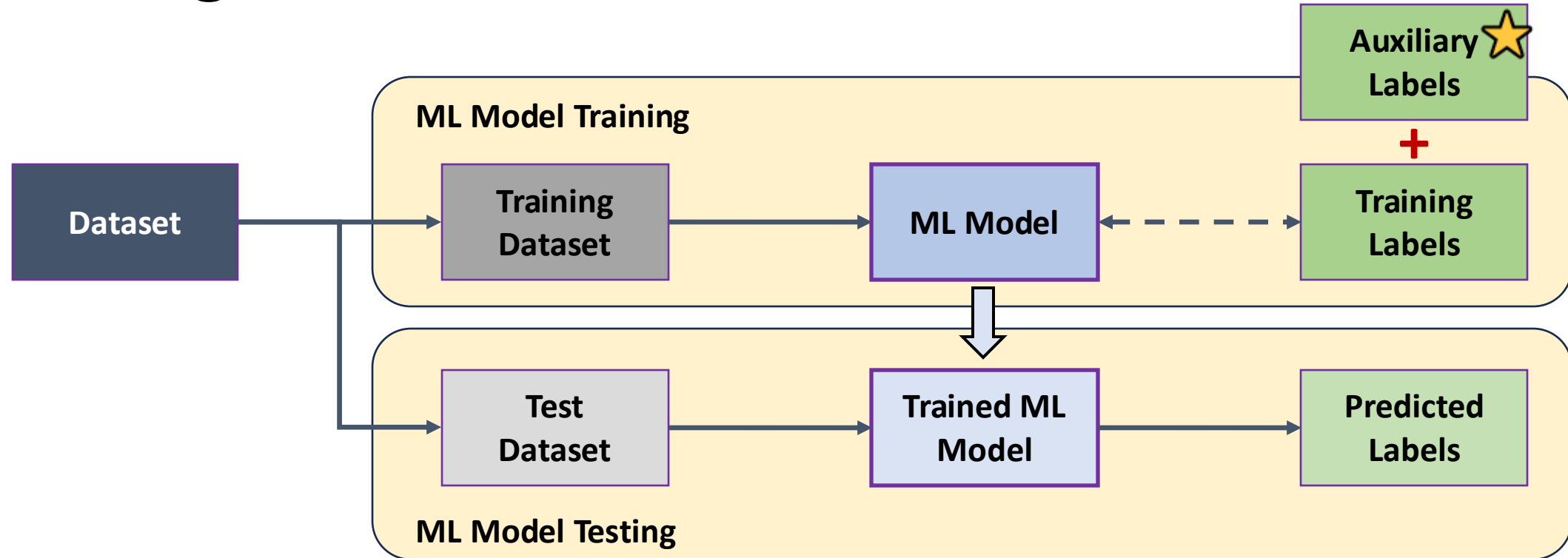


AARHUS UNIVERSITY



YOUNG INTERNATIONAL ACADEMICS

Background



- Introducing *auxiliary learning tasks* has been shown to be an effective way to address the problem of limited training labels in many applications.
- For example, incorporating chemical, physical, and toxicological profiles as auxiliary tasks helps ML models capture underlying data structures more effectively and improve generalization.

Motivation

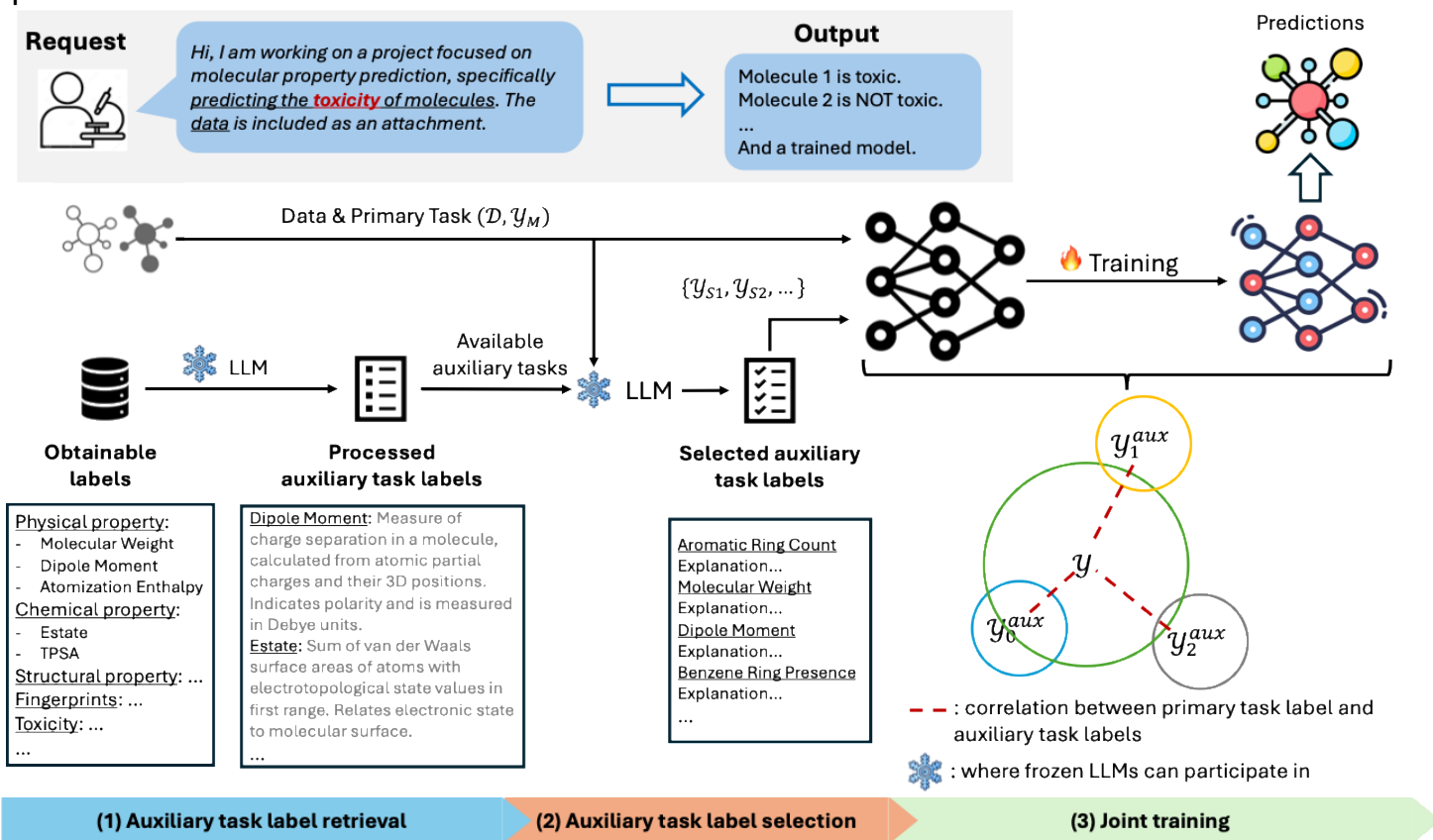
MODEL	AUTOMATIC TASK RETRIEVAL	AUTOMATIC TASK SELECTION	ADAPTIVE WEIGHTING	SELF-CONTAINED	NECESSARY INPUTS
UNWEIGHTED AVERAGES	✗	✗	✗	✓	AUXILIARY TASKS
TAG [8]	✗	✗	✓	✓	AUXILIARY TASKS
TASK2VEC [1]	✗	✗	✓	✓	AUXILIARY TASKS & DESCRIPTIONS
MTDNN [25]	✗	✗	✗	✗	RELEVANT AUXILIARY DATASETS & DESCRIPTIONS
GRADNORM [3]	✗	✗	✓	✓	AUXILIARY TASKS
GS-META [68]	✗	✓	✗	✗	AUXILIARY TASKS & RELATION GRAPH
MOLGROUP [17]	✗	✗	✓	✗	RELEVANT AUXILIARY DATASETS & DESCRIPTIONS
INSTRUCTMOL [51]	✗	✗	✓	✗	RELEVANT AUXILIARY DATASETS & DESCRIPTIONS
AUTAUT (OURS)	✓	✓	✓	✓	NONE

- However, constructing high-quality auxiliary tasks is often complex and resource-intensive.
- In fields such as biology and chemistry, these challenges are further amplified by the scarcity, high cost, and time demands of obtaining domain-specific knowledge.

RQ: Can we automate the process of auxiliary task-enhanced ML model training?

Method - AutAuT

TL;DR: A fully automated framework that uses LLMs to retrieve *Auxiliary Tasks* and employs *Adaptive Weighting* to integrate them for molecular property prediction.



1) Auxiliary Task Retrieval

- 1) *Instruction*: Provides general guidance to the LLM, specifying its role in the retrieval process.
- 2) *Message*: A direct and clear request for the LLM to identify potential auxiliary task labels based on the given context

2) Auxiliary Task Selection

- 1) Search available information about these properties and write a brief summary...
- 2) Assess the relevance between retrieved properties and the primary task...
- 3) For the primary task, recommend K properties as auxiliary tasks...

Method - AutAuT

Name	Category	Brief Description
Molecular Weight	Constitutional	Total mass of a molecule calculated as the sum of atomic weights of all atoms, providing fundamental information about molecular size and mass distribution.
Heavy Atom Molecular Weight	Constitutional	Sum of atomic weights of all non-hydrogen atoms in the molecule, useful for comparing core molecular frameworks.
Number of Valence Electrons	Constitutional	Total number of electrons in the outer shells of all atoms, crucial for understanding chemical bonding and reactivity patterns.
Total Formal Charge	Constitutional	Sum of all formal charges on atoms in the molecule, indicating overall molecular charge state and ionic character.
Topological Polar Surface Area	Topological	Sum of surfaces of all polar atoms (mainly oxygen and nitrogen), correlating with drug absorption, including intestinal absorption and blood-brain barrier penetration.
Labute Approximate Surface Area	Topological	Approximate molecular surface area calculated using Labute's method, useful for predicting physical properties and molecular interactions.
Balaban J Index	Topological	Topological index based on molecular connectivity, indicating molecular branching and cyclicality. Higher values suggest more branched structures.
Bertz Complexity	Topological	Measure of molecular complexity considering both size and branching patterns. Higher values indicate more complex molecular structures.
LogP	Electronic	Logarithm of octanol-water partition coefficient, predicting molecular lipophilicity and membrane permeability. Key for drug absorption.
Molar Refractivity	Electronic	Measure of total polarizability of a molecule, related to molecular volume and electronic properties. Important for predicting optical behavior.
EState VSA1	Electronic	Sum of van der Waals surface areas of atoms with electrotopological state values in first range. Relates electronic state to molecular surface.

1) Auxiliary Task Retrieval

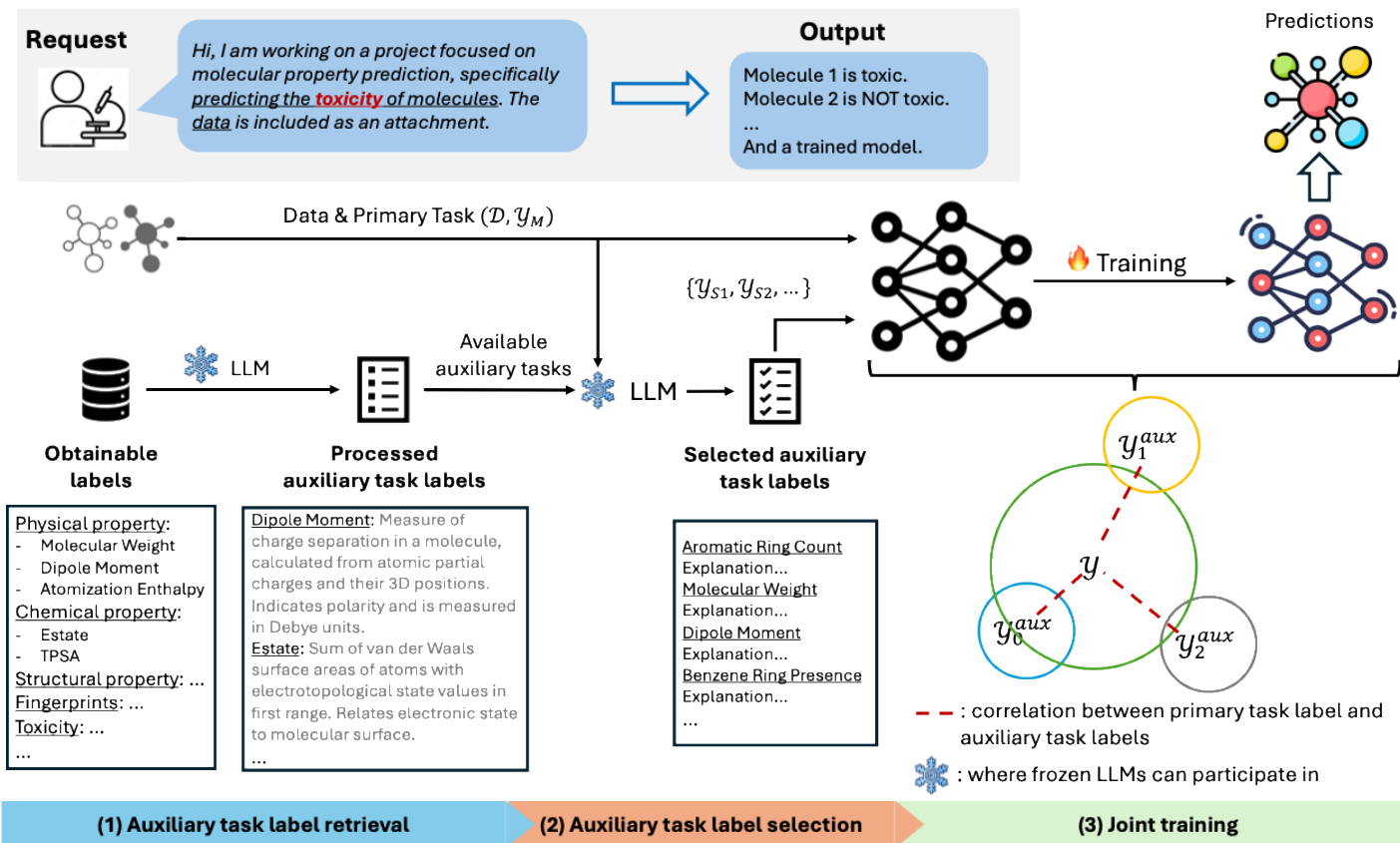
- 1) *Instruction*: Provides general guidance to the LLM, specifying its role in the retrieval process.
- 2) *Message*: A direct and clear request for the LLM to identify potential auxiliary task labels based on the given context

2) Auxiliary Task Selection

- 1) Search available information about these properties and write a brief summary...
- 2) Assess the relevance between retrieved properties and the primary task...
- 3) For the primary task, recommend K properties as auxiliary tasks...

Method - AutAuT

TL;DR: A fully automated framework that uses LLMs to retrieve *Auxiliary Tasks* and employs *Adaptive Weighting* to integrate them for molecular property prediction.



- 1) Auxiliary Task Retrieval
- 2) Auxiliary Task Selection
- 3) Learning to Learn from Selected Auxiliary Tasks

- 1) *Weight initialization*: the auxiliary task weights are initialized based on affinity scores.
- 2) *Dynamic weight adaptation*: the auxiliary task weights are iteratively updated using gradient alignment and validation performance.
- 3) *Dynamic weight adaptation*: the ML model parameters are optimized with fixed auxiliary task weights.

Theorem 1

$$\|\nabla \mathcal{L}_M(\theta) - \sum_{k=1}^K \alpha_k \nabla \mathcal{L}_S^k(\theta)\| = 0.$$

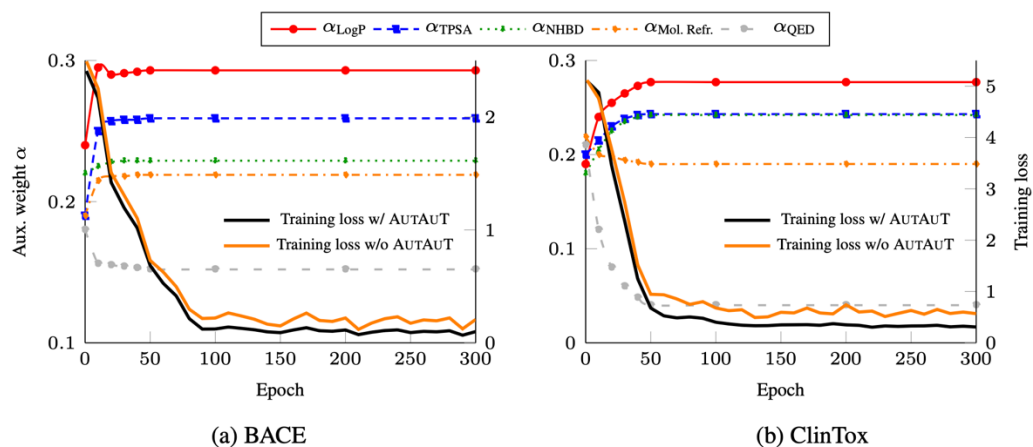
(Subspace Alignment for Optimization)

Theorem 2

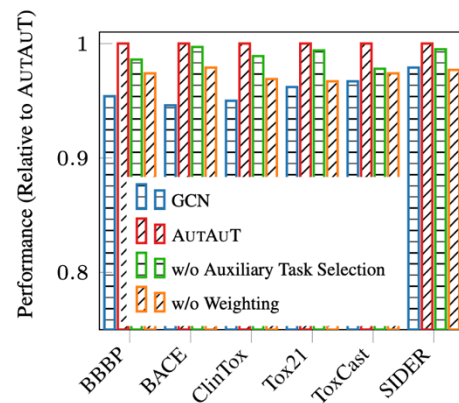
$$\mathcal{E}(\theta) \leq \hat{\mathcal{E}}(\theta) + c \cdot \mathcal{R}_n(\mathcal{H}_\alpha),$$

(Generalization Bounds with Auxiliary Tasks)

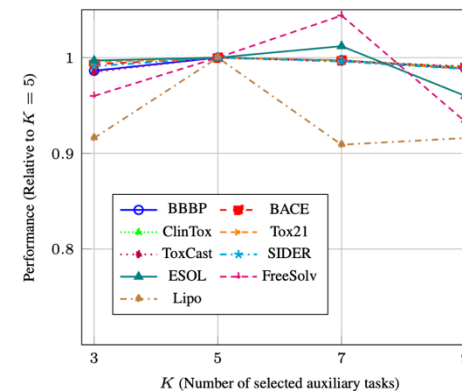
Experiments



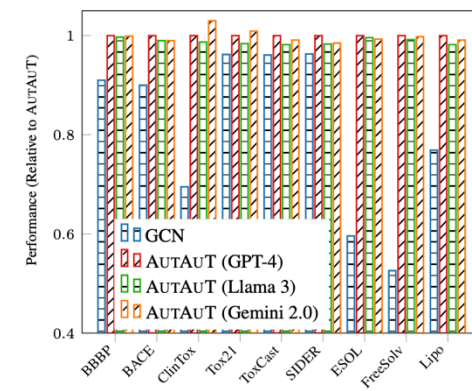
A



B



C



D

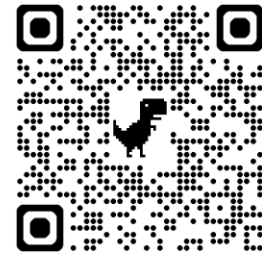
A: Auxiliary task influence is correctly inferred.

B: Auxiliary task selection and adaptive weighting complement each other.

C: The number of selected auxiliary tasks (K) balance informativeness and noise.

D: All major LLMs produce useful task suggestions and improve over baseline.

Thank you!



Contact



Paper