Ron

# Some papers in chemistry and machine learning

30 June 2025

# Index

- Tokenization
- NovelSeek
- Tools
- ChemOCR

# Tokenization

A way to feed molecules to machines

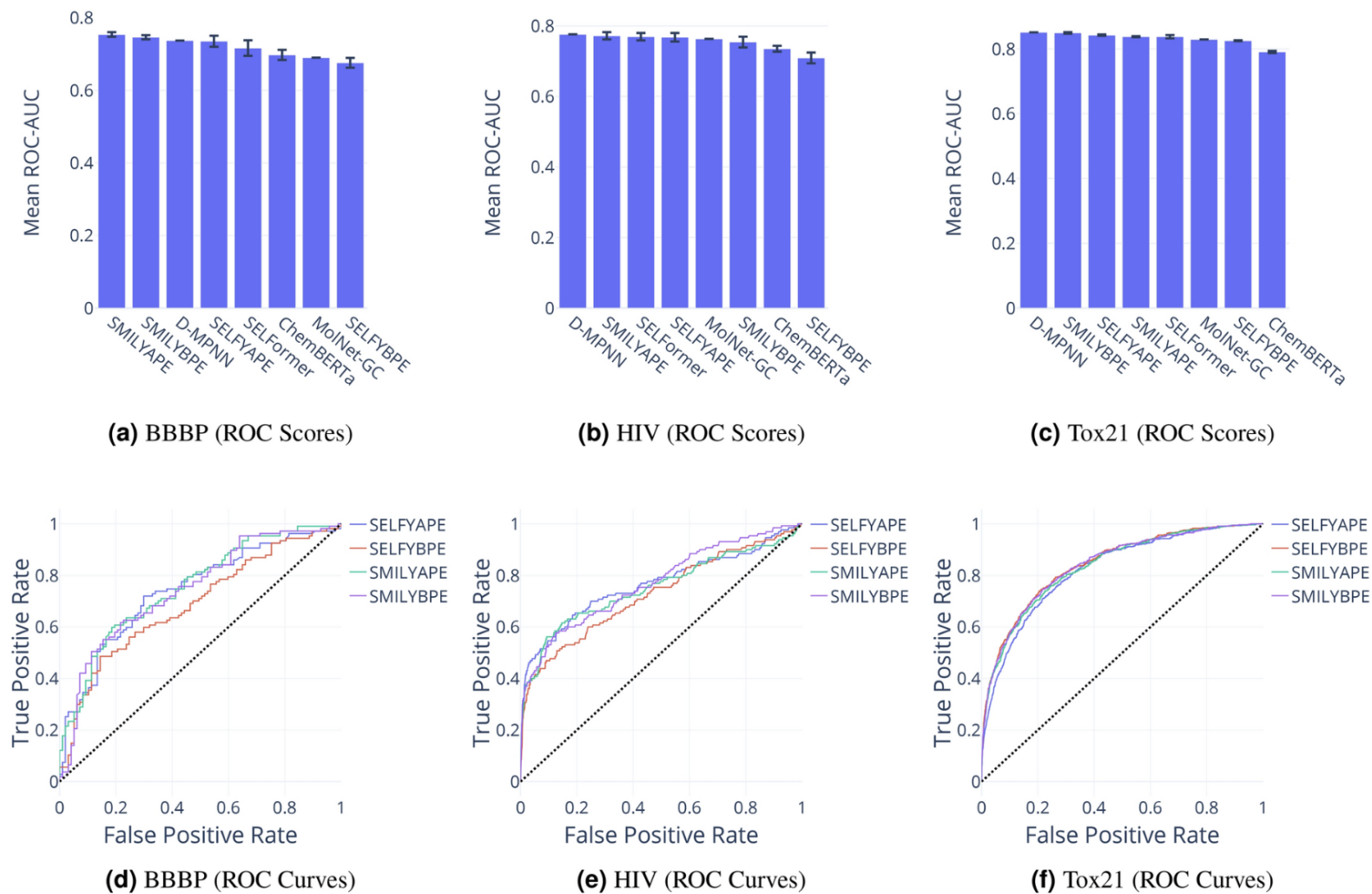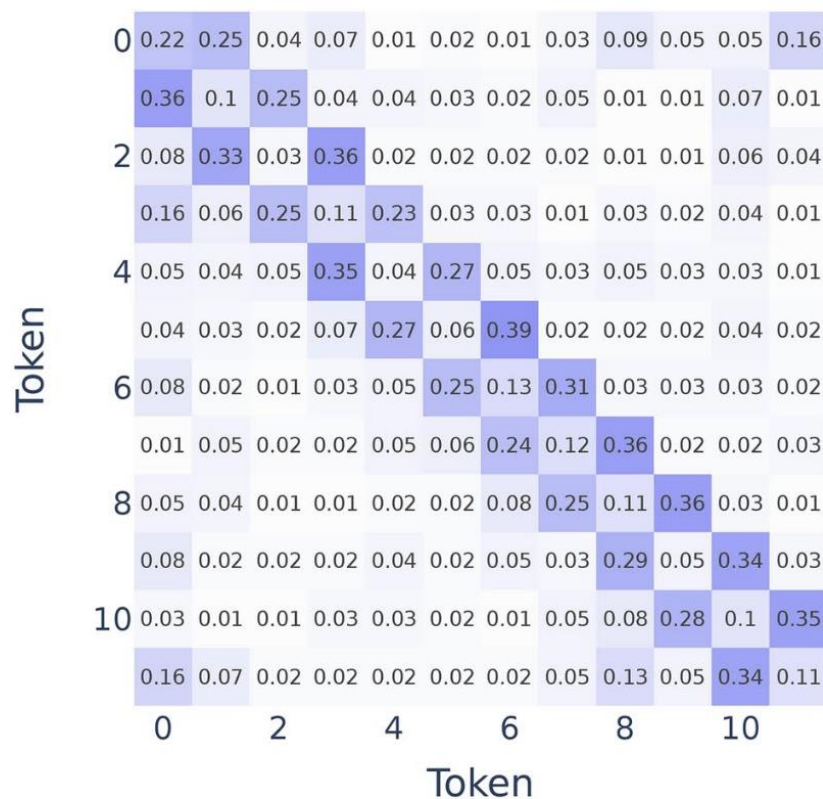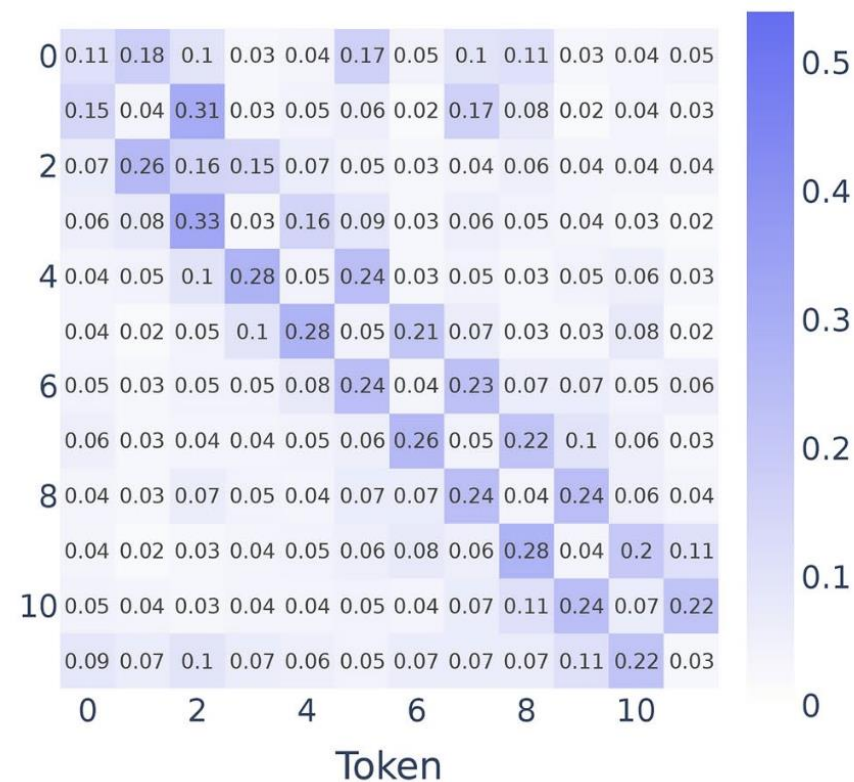# Is there a difference in performance?

(Not much)

No.

**Fig. 5**. ROC-AUC scores and curves of selected benchmarks. Top row: (**a**) BBBP, (**b**) HIV, and (**c**) Tox21 bar charts. Bottom row: (**d**) BBBP, (**e**) HIV, and (**f**) Tox21. The Tox21 is the largest dataset, giving better results with all models.
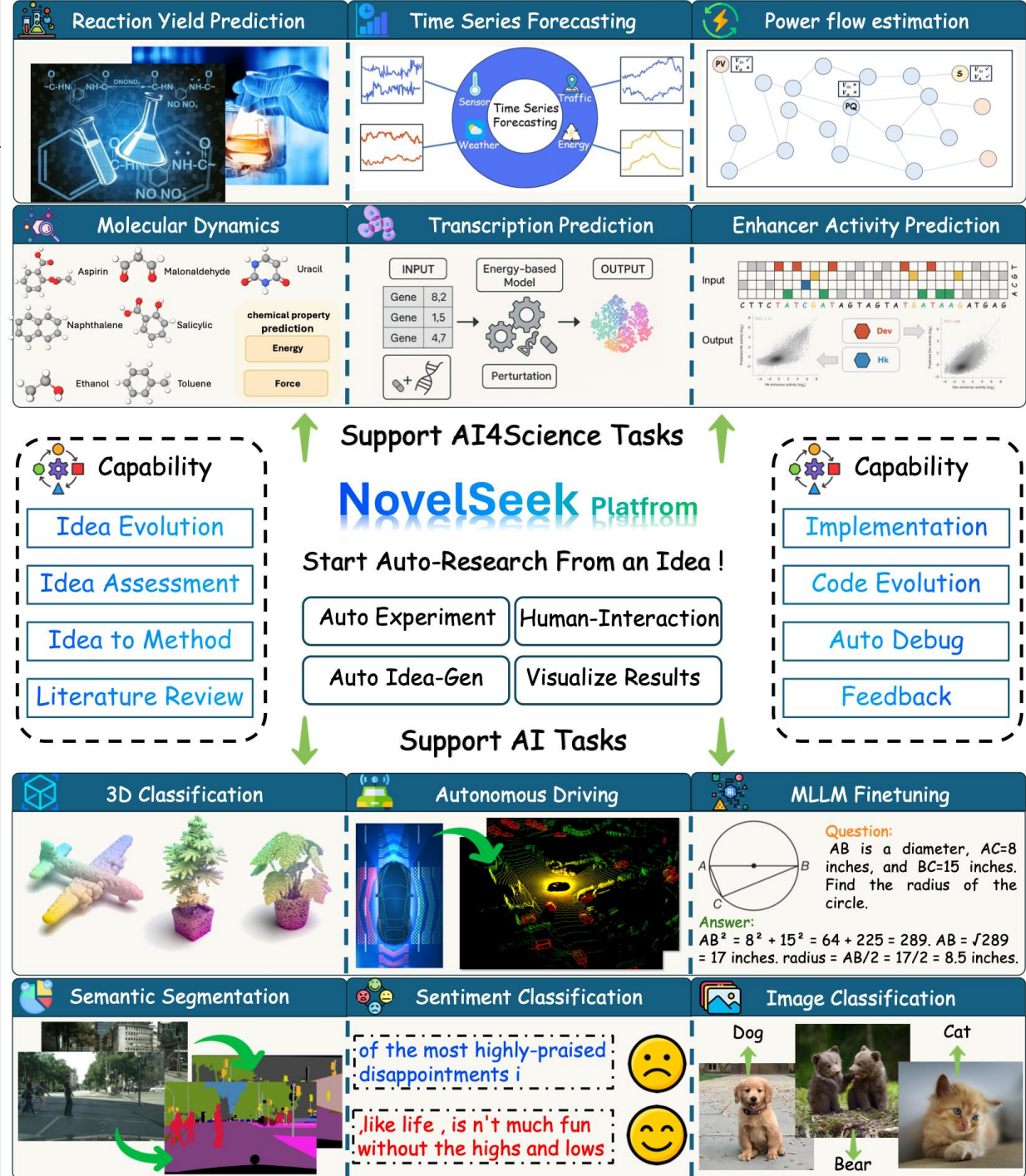
**(a)** SMILYAPE $C_{17}H_{21}NO_2$

**(b)** SELFYAPE $C_{17}H_{21}NO_2$

**Fig. 6**. Attention patterns of the SMILYAPE (**a**) and SELFYAPE (**b**) models, respectively, for molecule $C_{17}H_{21}NO_2$ from the BBBP dataset. The visualization highlights a stronger focus on attention between adjacent tokens (darker squares along the diagonal). However, some attention is also directed towards distant tokens (lighter squares), suggesting the models capture both local context and broader semantic relationships within the molecule.

# NovelSeek

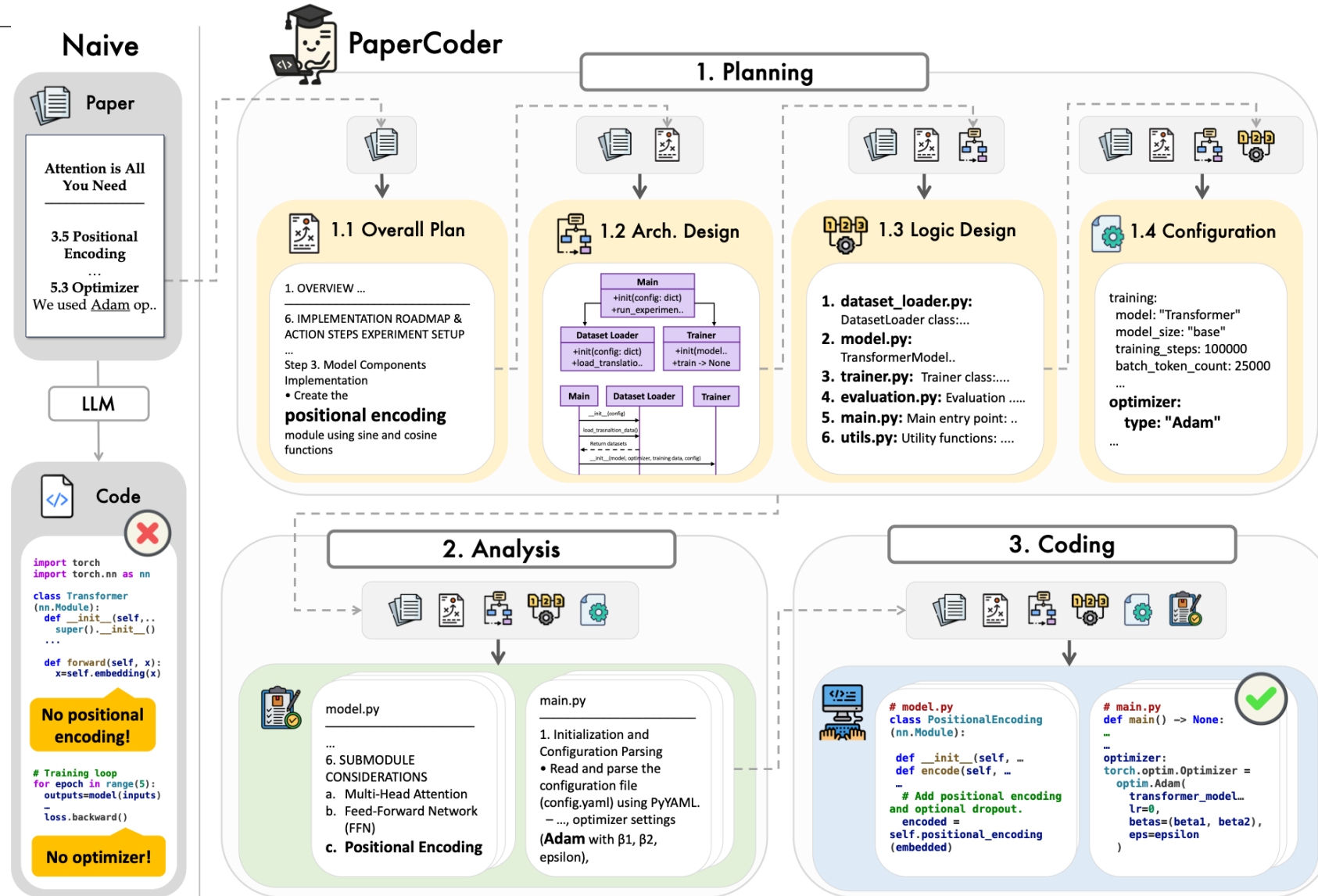When Agent Becomes the Scientist – Building Closed-Loop System from Hypothesis to Verification
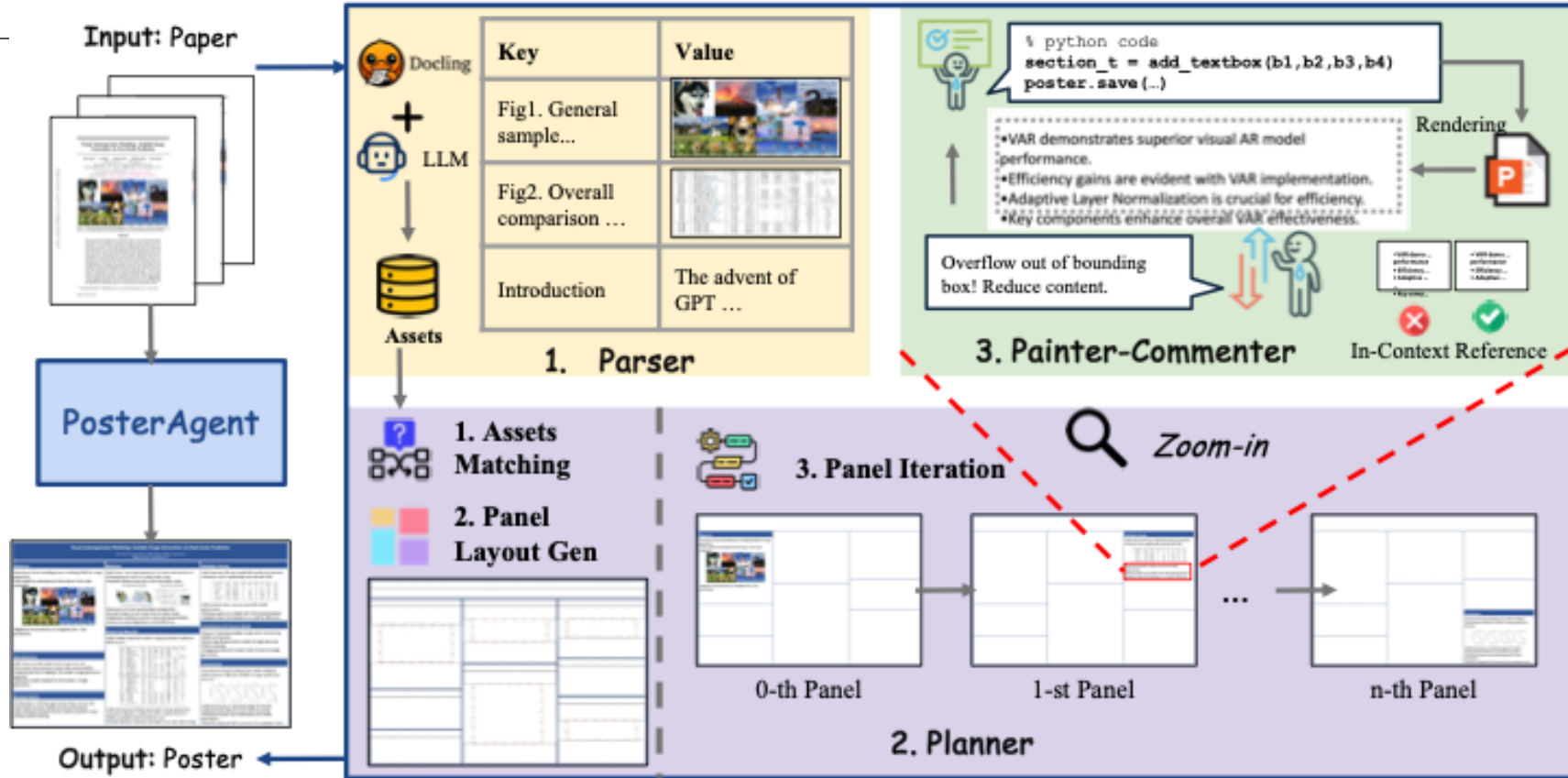
# Some useful tools

(For you?)

# Paper2Code

Replicates code from a(n ML) paper autonomously.

# Paper2Poster

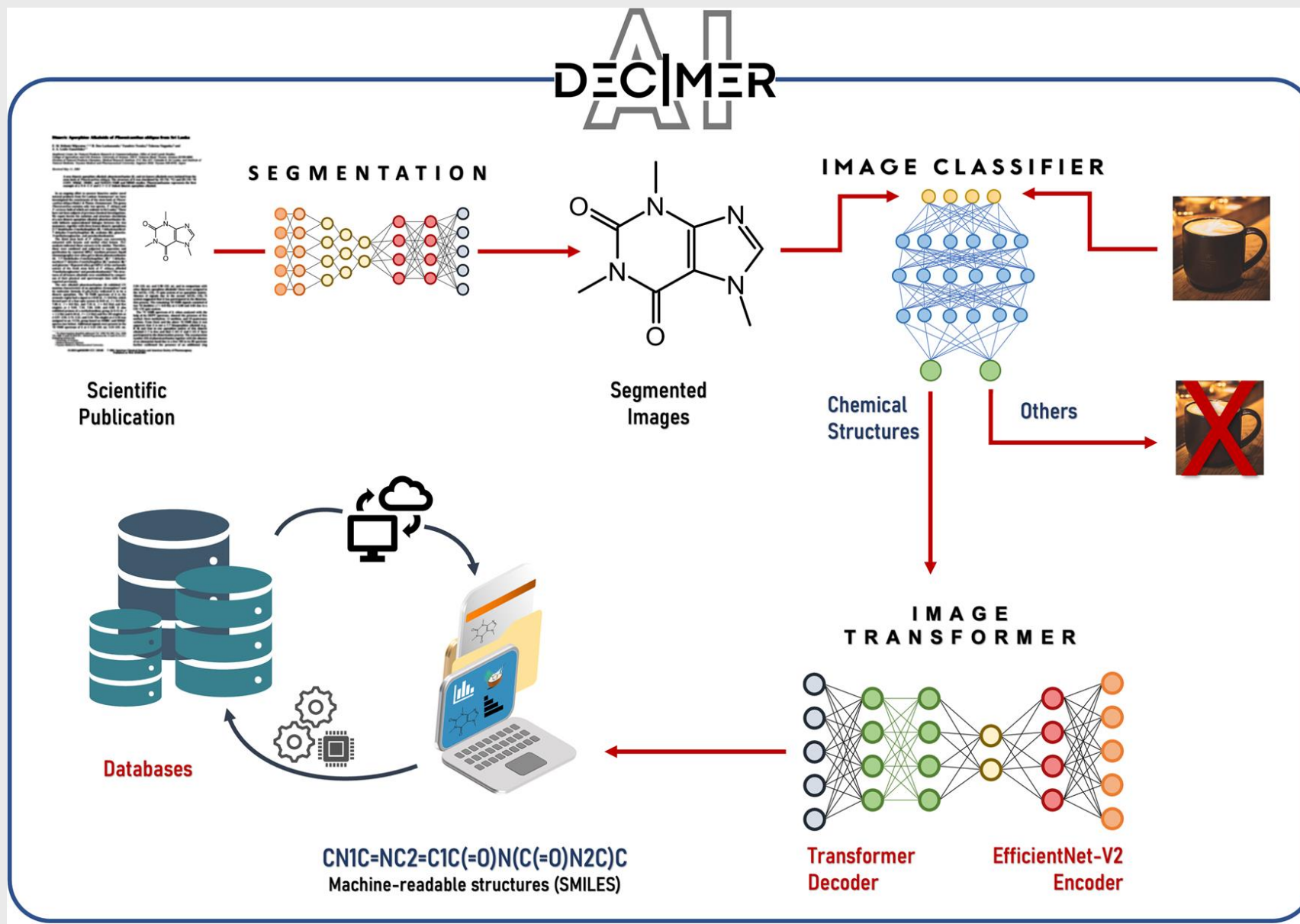Converts papers to poster autonomously (and efficiently).

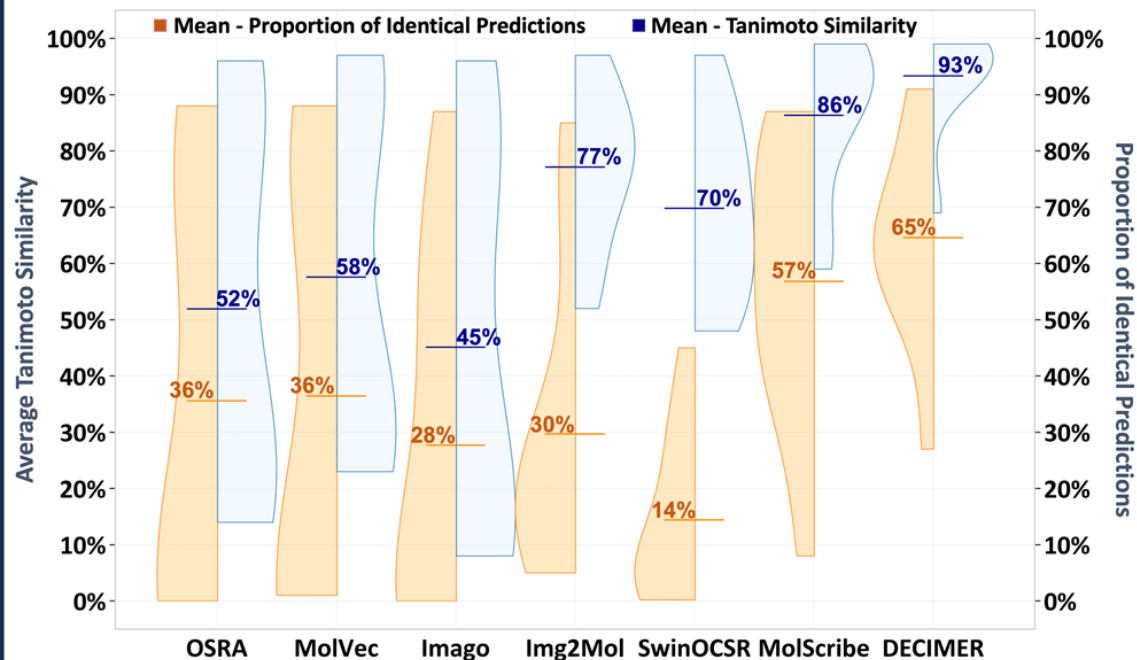# Identifying molecules from images

## Motivation

- As mentioned in a previous journal club presentation, a large amount of chemical knowledge is documented on paper or scanned images of them.
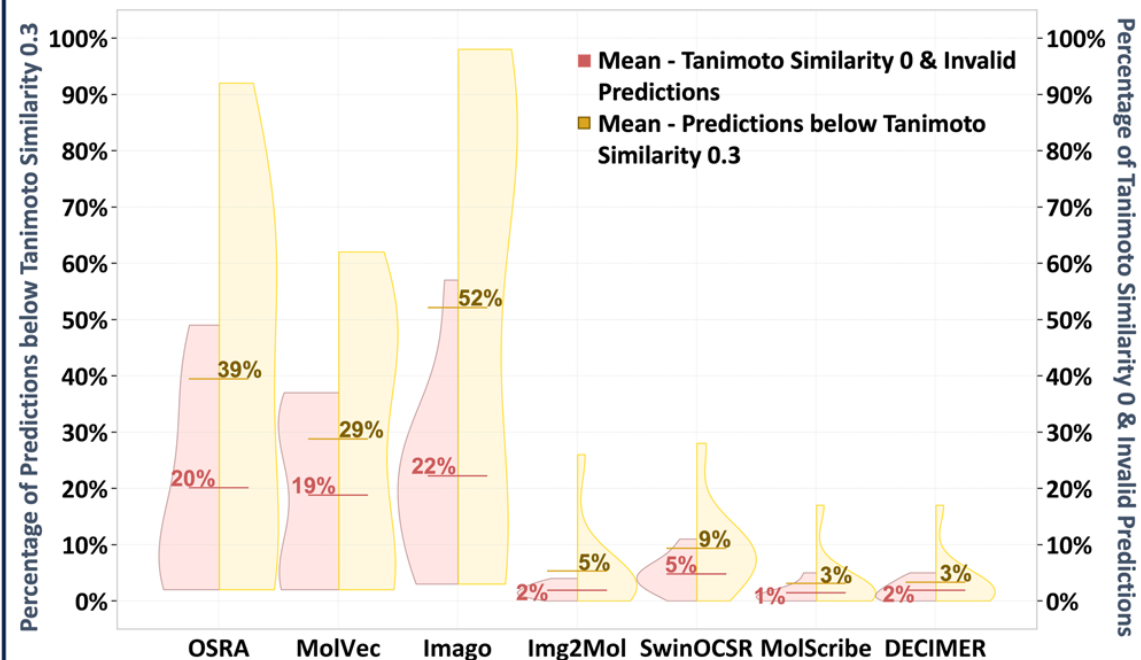- Efficiently digitizing them would help models absorb their knowledge.
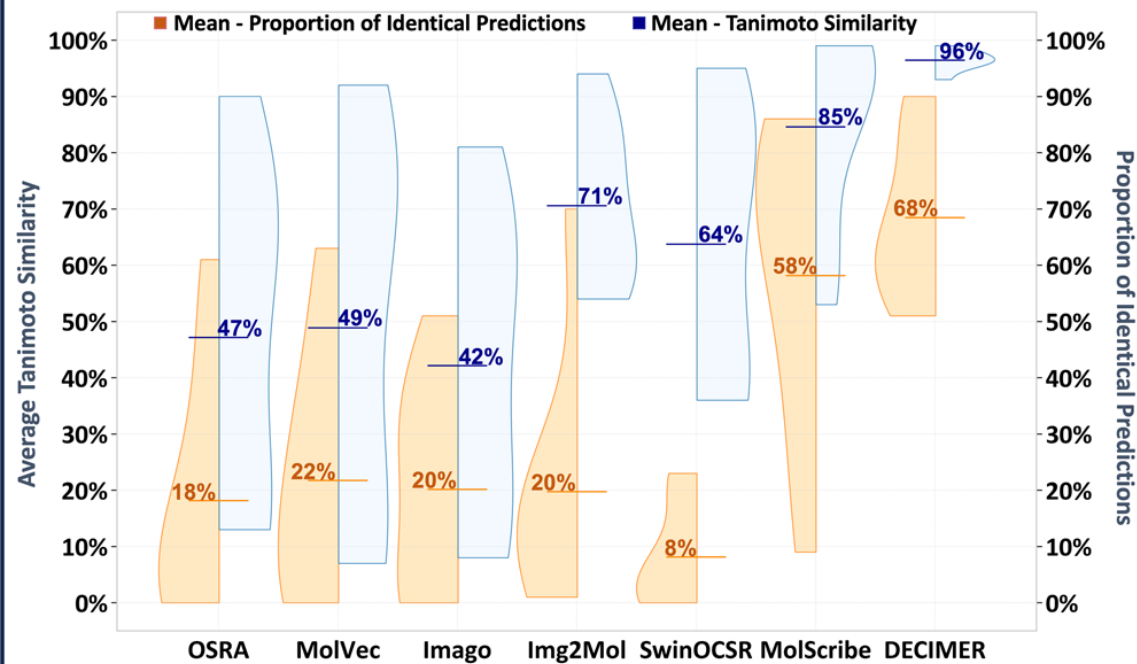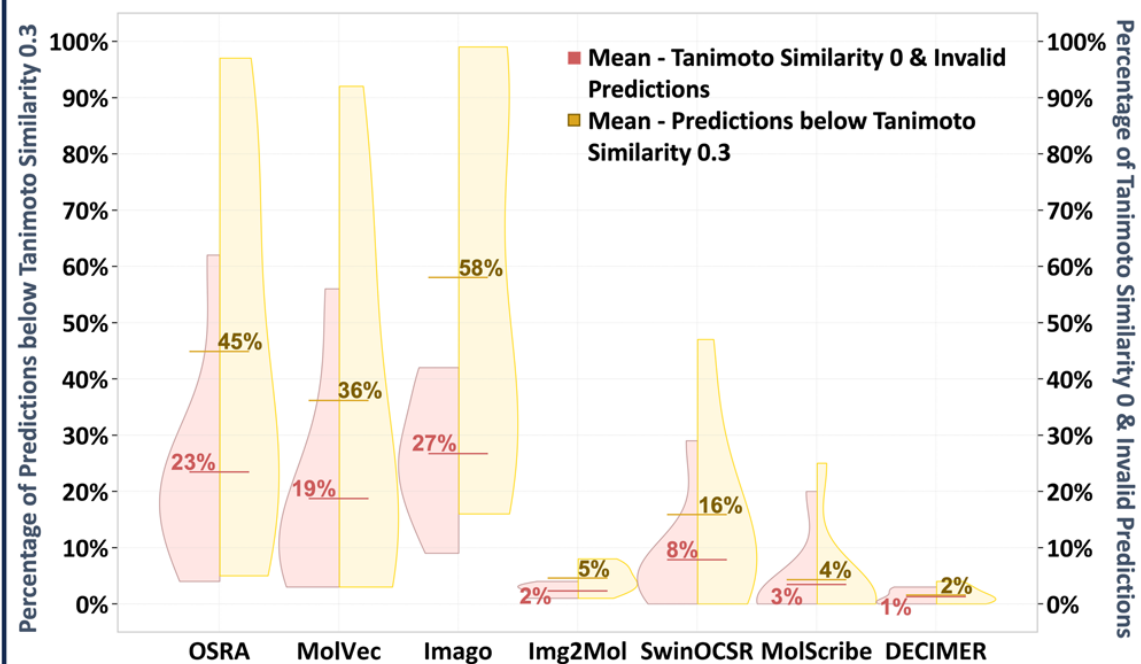
A: Success rates for datasets without added distortions

B: Success rates for datasets with added distortions

C: Failure rates for datasets without added distortions

D: Failure rates for datasets with added distortions

# Thank You