

Unveiling Hidden Treasures: Exploring Large Chemical Spaces with Machine Learning Models Trained on DNA-Encoded Libraries Selection Results

A. Kapeliukha¹, O. Strapak¹, O. Tarkhanova¹, S. Zozulya³, Alla Pohribna³, Yurii S. Moroz²

¹Chemspace LLC, 85 Winston Churchill Street, Kyiv 02094, Ukraine

²Enamine Ltd, 78 Winston Churchill Street, Kyiv 02094, Ukraine

³Bienta, 78 Winston Churchill Street, Kyiv 02094, Ukraine

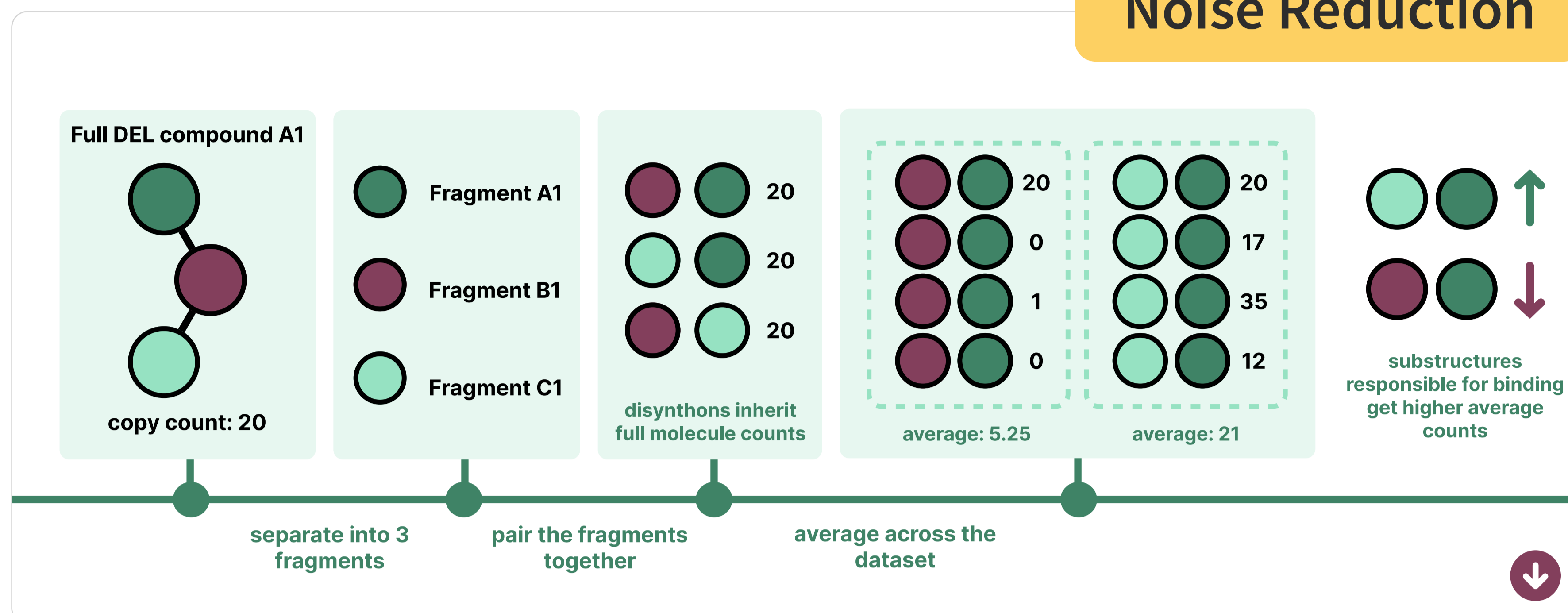
DEL Screening Data

DEL screening results:

- ✓ 4 experiment selections
- ✓ 2 control selections
- ✓ 108.5k datapoints

CAIX

Noise Reduction



Regression Model Performance

Full molecules - prediction on validation set

Spearman's r: 0.28 ± 0.01

Disynths - prediction on validation set

Spearman's r: 0.59 ± 0.01

Calculated enrichment

Calculated enrichment

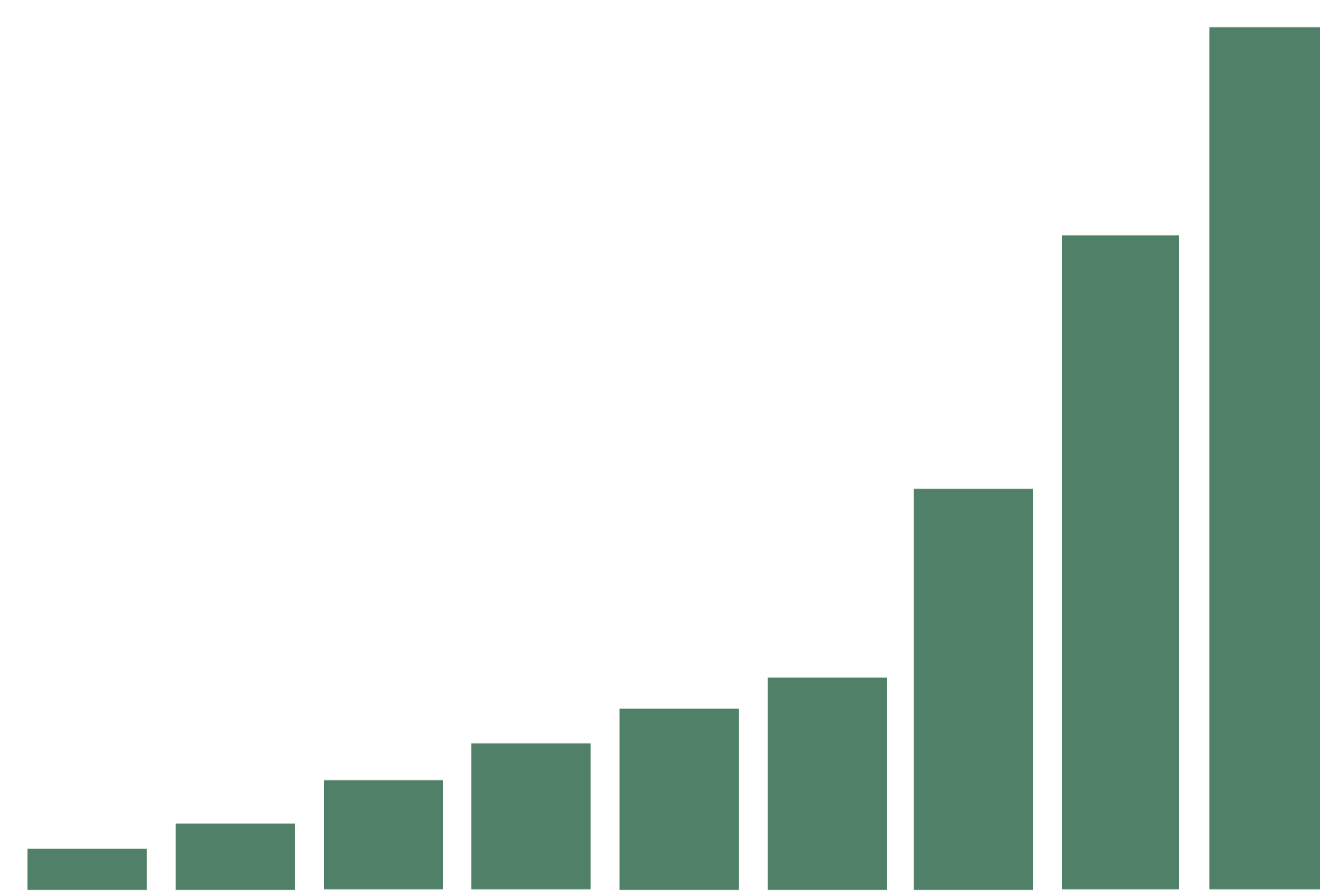
Predicted enrichment

Predicted enrichment

Selecting active compounds from ChEMBL

Number of active compounds in top N ChEMBL ranked by predictions

N actives identified



Top N ranked ChEMBL compounds

Exploration of Enamine REAL

33B+
Enamine REAL
Space

1M
Top compounds by
prediction

Sphere
exclusion
clustering

Visual
inspection

- ✓ Regression model with best performance used for prediction
- ✓ Selected diverse compounds with the best predicted enrichment in the cluster
- ✓ Compounds are dissimilar from known actives
- ✓ Synthesis success rate 87% by Enamine

150
Selected for
synthesis

130
Successfully
synthesized

Results

- ✓ Disynthon aggregation helps to identify substructures important for binding and de-noise DEL datasets
- ✓ Regression models trained on DEL data can efficiently rank compounds from large compound databases.
- ✓ Exploration on Enamine REAL by ML model leads to identification of 33 novel hits with 25.4% hit rate.

IC50 = 0,52 μM

IC50 = 0,45 μM

25,4%
hit rate

IC50 = 0,13 μM

IC50 = 0,28 μM

a.kapeliukha@chem-space.com

www.chem-space.com