




NeurIPS 2024 - Predict New Medicines with BELKA


Predict small molecule-protein interactions using the Big Encoded Library for Chemical Assessment

Background

 20+ years of corporate finance experience. BS in economics, MS in management.

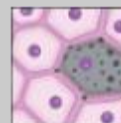
 Ex-CFO of a Russian version of Fannie Mae (\$20bn assets).

 MSc candidate in Data Science and AI at Northwestern (2025).

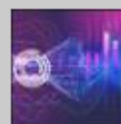
 Primary focus: medical imaging, chemistry, pharmaceuticals, and [just on the edge of the radar] finance.



Predict small molecule-protein interactions using the Big Encoded Library for Chemical Assessment: **1st out of 1950 teams**



Hacking the Human Vasculature in 3D. Segment vasculature in 3D scans of human kidney: **18th out of 1149 teams**



Jane Street Real-Time Market Data Forecasting. Predict financial market responders using real-world data: **7th out of 2280 teams****

Challenge: a needle in a haystack

C#CCCC[C@H](Nc1nc(NCCCNC(=O)c2occc2C)nc(Nc2sc(C1)c2C(=O)OC)n1)C(=O)N[Dy]

C#CCCC[C@H](Nc1nc(NCCN2C(=O)Cc3ccccc32)nc(Nc2c(C)cc([N+](=O)[O-])cc2OC)n1)C(=O)N[Dy]

C#CCCC[C@H](Nc1nc(NC2cc(C1)s2C1)nc(Nc2cc(C1)nc(SC)n2)n1)C(=O)N[Dy]

C#CCCC[C@H](Nc1nc(NC2nnc(S)o2)nc(Nc2ccc(C=C)cc2)n1)C(=O)N[Dy]

C#CCCC[C@H](Nc1nc(NCCC2(COC)CC2)nc(Nc2ccc(N3CC3)cc2)n1)C(=O)N[Dy]

C#CCCC[C@H](Nc1nc(NCCc2nc[nH]c2C)nc(Nc2cccc3c2C(=O)NC3=O)n1)C(=O)N[Dy]

C#CCCC[C@H](Nc1nc(NCC2CCC(SC)CC2)nc(Nc2ncnn3cccc23)n1)C(=O)N[Dy]

C#CCCC[C@H](Nc1nc(NC2ccn3ccnc3c2)nc(NC(=N)c2cccc(CN)c2)n1)C(=O)N[Dy]

C#CCCC[C@H](Nc1nc(NCCc2csc3ccc23)nc(NC[C@H]2CC[C@H](c3cccc3)O2)n1)C(=O)N[Dy]

C#CCCC[C@H](Nc1nc(NC2ccns2)nc(Nc2sc(C1)cc2C(=O)OC)n1)C(=O)N[Dy]

C#CCCC[C@H](Nc1nc(NCCc2ccc(S(N)(=O)=O)s2)nc(NC2cc(C1)sc2C1)n1)C(=O)N[Dy]

C#CCCC[C@H](Nc1nc(NC2ccncc2)nc(NC2nnc(N3CC3)o2)n1)C(=O)N[Dy]

C#CCCC[C@H](Nc1nc(Nc2ccc3[nH]cnc3c2)nc(Nc2cc(C(F)(F)F)ccn2)n1)C(=O)N[Dy]



Drug-like-space: **10⁶⁰** molecules

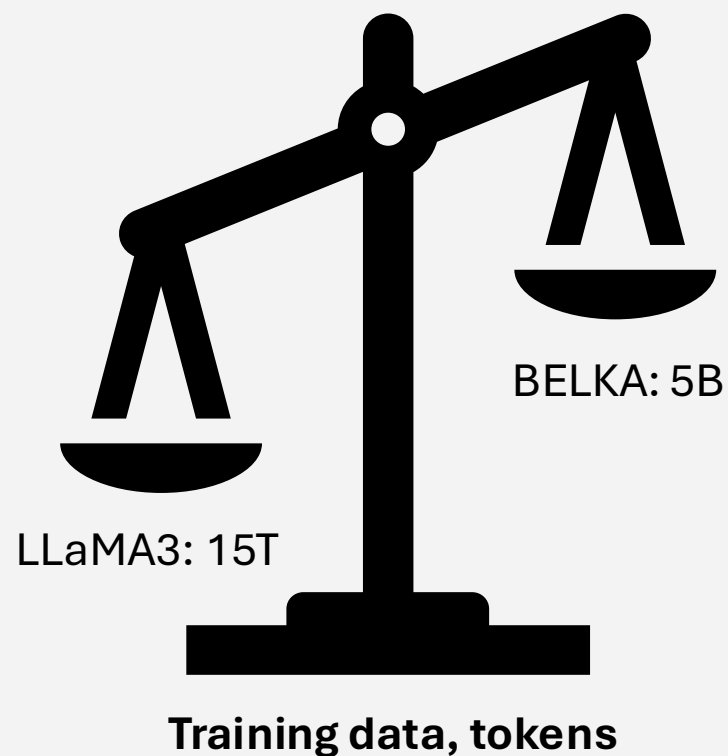


BELKA: **133M** small molecule drug candidates, 3 target proteins



FDA-registered drugs: **4471**

NLP vs. Chemistry: 3000-0 to Poets...



	NLP	Chem
Vocabulary	170K words, 50K+ tokens	118 periodic elements
Rules / structure	Vague and evolving	Defined and universal
Targets	Vague and context-driven	Clear
Uncertainty	High: polysemy, ambiguity, cultural context	Moderate: fundamentally deterministic

Architecture [finally]

Chem-aware tokenization.

43 tokens: Br, C, [C++], etc.
Also tried but failed: n-grams,
Wordpiece, Bert.

Embedding / latent space

dimension: 32

ChemBERTa-2²: 591 tokens,
embedding size 768.

1) <https://pypi.org/project/atomInSmiles/>

2) <https://arxiv.org/abs/2209.01712>



Tokenizer: atomInSmiles¹

Embeddings /
Positional Encodings

4 x Encoder (Self-
Attention -> FFN -> Add)

Neck: GRU

Head: Dense(sigmoid)

Training schedule

MLM: classics from “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding”.
Mini-batch adaptive Categorical Focal Cross Entropy.

ECFP and Affinity: GRU feature extractor on top of Self-Attention Encoder.
Binary Focal Cross Entropy

	MLM	ECFP	Affinity
Train data	✓	✓	✓
Public test data	✓	✓	
Extra data ¹	✓	✓	
Epochs (100M each)	10+	5+	10+
Validation	None	9m non-shared	

1) "Building Block-Based Binding Predictions for DNA-Encoded Libraries", cited by @hengck23 and processed by @chemdatafarmer







Brilliant ideas [that failed]

**Pretrain to predict
outcomes of chemical
reactions**

**GAN-based
augmentations**

**SMILES-to 3D
fingerprints pre-training**

**Read "Chemistry for
Dummies"**

-  ZINC pre-training
-  MTR pre-training
-  Fixed-position fingerprints (MACCS, PubChem)
-  ECFP- or multi-input models
-  High-dimensional models
-  Gated Fusion / Cross-attention with blocks

Contacts:

vshlepov@gmail.com

<https://www.linkedin.com/in/victor-shlepov/>