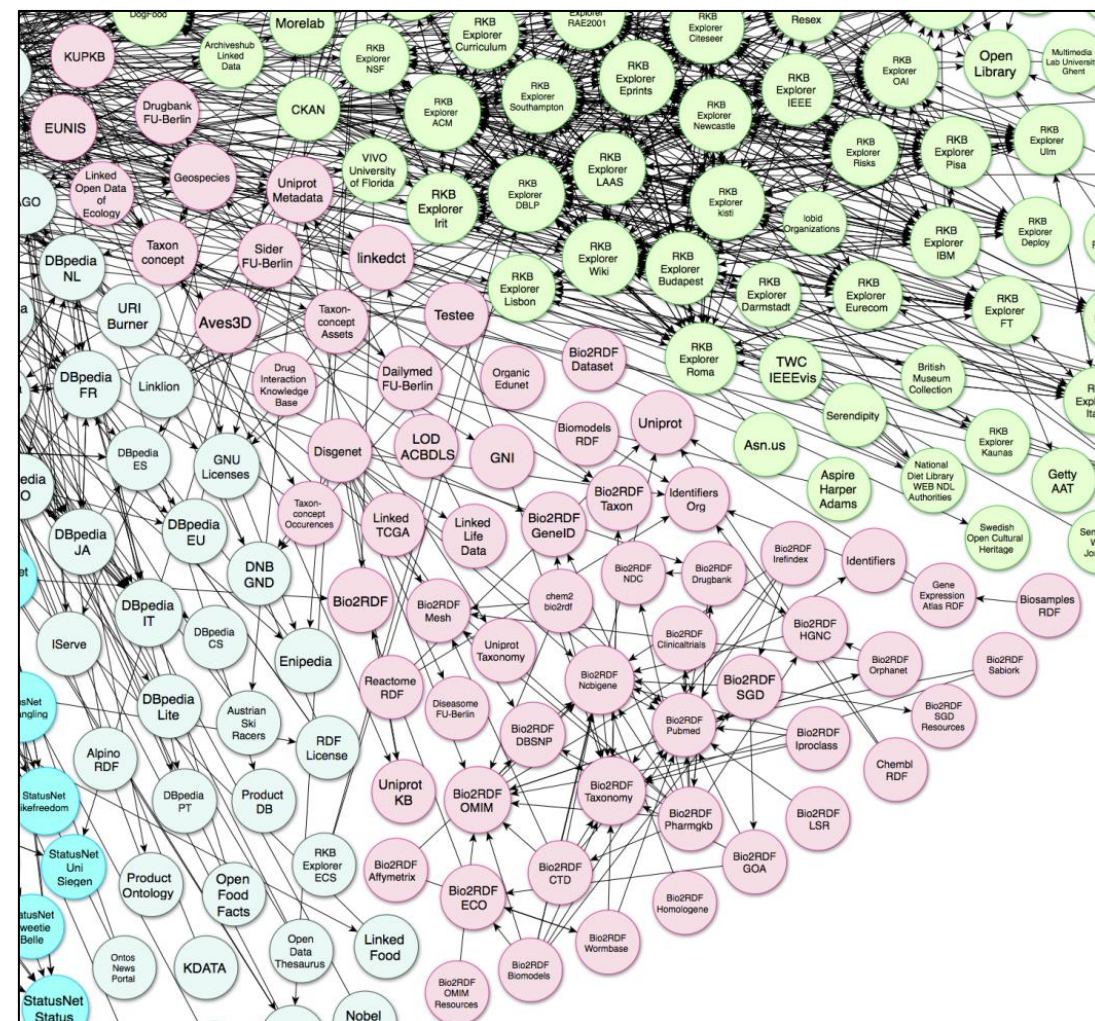## Motivation

- Questions asked during the Hypothesize-test-evaluate cycle: *"Which antineoplastic agents target IDH1 gene in glioma patients?"*
- Requires precise answers (in above example, the corresponding drugs) or relevant *-omics* datasets for further analysis.
- Semantic Web and Linked Data technologies (RDF and SPARQL) towards tackling the integrative bioinformatics challenges.
- Querying Life Sciences Linked Open Data (LSLOD) Cloud has a steep learning curve.
- Natural language querying method over the LSLOD network to enable scalable, autonomous discovery of relevant answers and datasets for evaluating hypotheses.
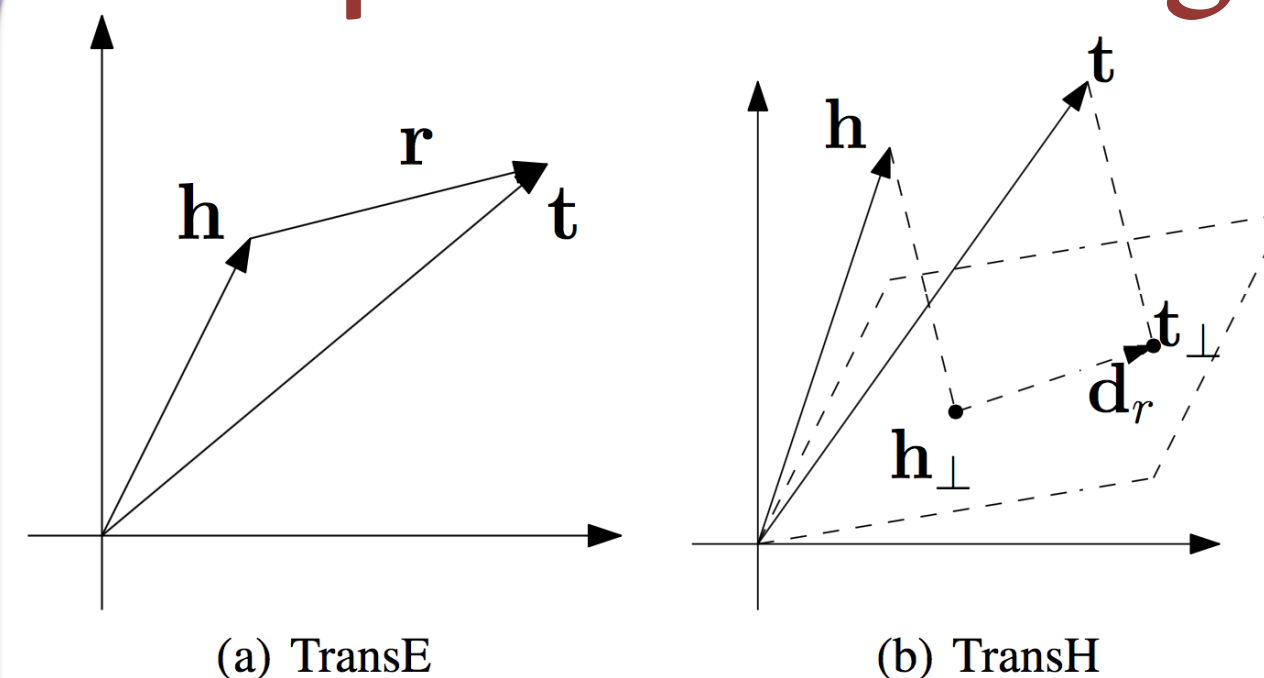
## Word Embeddings

**PubMed**
**MEDLINE** U.S. National Library of Medicine

1.5B+ abstracts, 3B+ Tokens
2.2M+ unique words
0.14M biomedical entities

**PubTator** → **Annotation**



**↓ Tokenization**

**Effect of hemodialysis on methylprednisolone plasma levels.** The effect of hemodialysis on methylprednisolone levels in uremia was investigated. Methylprednisolone 15 mg/kg was given intravenously over a period of 20 min ....

effect of hemodialysis on *id_chemical_d008775* plasma levels . the effect of hemodialysis on *id_chemical_d008775* levels in *id_disease_d014511* was investigated . *id_chemical_d008775* NNNNNNN mg/kg was given intravenously over a period of NNNNNNN min ....

## GloVe

Iterations: 20, $\alpha$: 0.75
Dimensions: 100

## Linked Data



Reference: http://lod-cloud.net/

- **Life Sciences Linked Open Data Cloud:** 1T+ triples from 80+ biomedical sources.
- Develop Bio-mashups and Linked Biomedical Dataspaces facilitating *in silico* data discovery
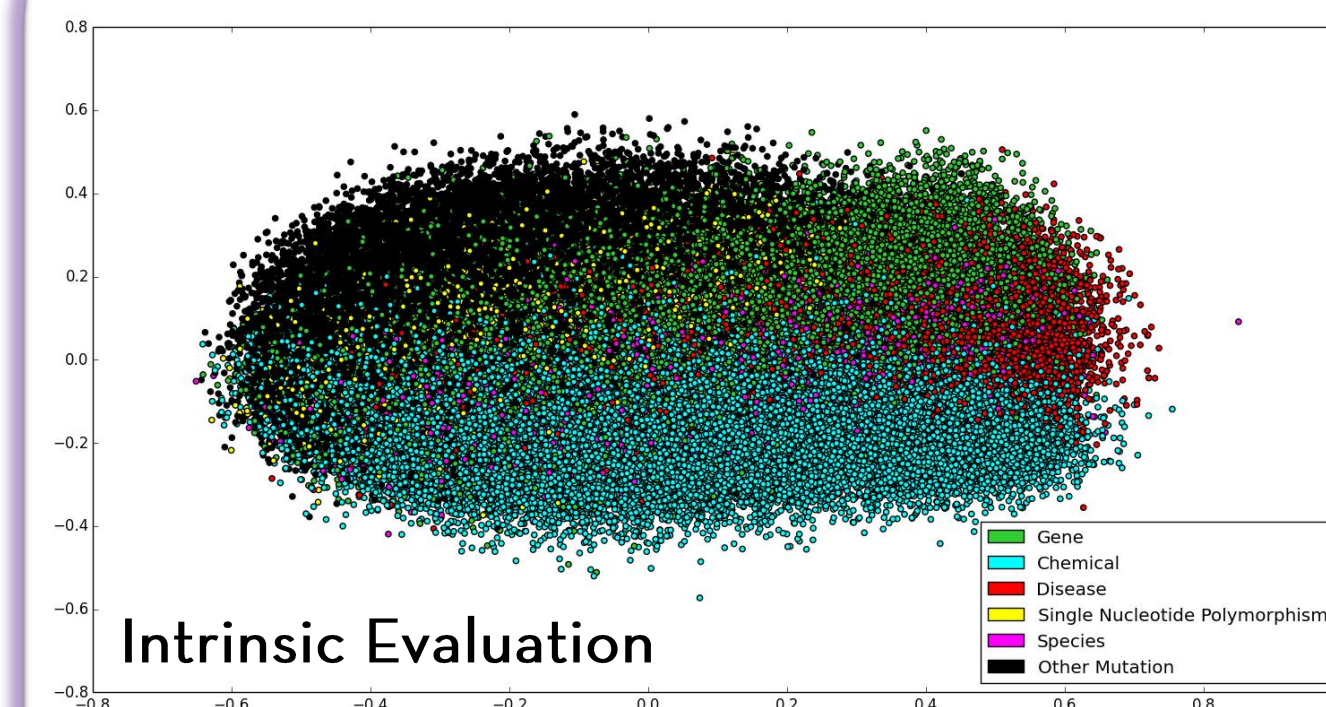
## Graph Embeddings



(a) TransE         (b) TransH

$$f_r(\mathbf{h}, \mathbf{t}) = \|(\mathbf{h} - \mathbf{w}_r^\top \mathbf{h} \mathbf{w}_r) + \mathbf{d}_r - (\mathbf{t} - \mathbf{w}_r^\top \mathbf{t} \mathbf{w}_r)\|_2^2$$

$$\mathcal{L} = \sum_{(h,r,t)\in\Delta} \sum_{(h',r',t')\in\Delta'_{(h,r,t)}} \left[ f_r(\mathbf{h}, \mathbf{t}) + \gamma - f_{r'}(\mathbf{h}', \mathbf{t}') \right]_+$$

$$+ C \left\{ \sum_{e\in E} \left[ \|\mathbf{e}\|_2^2 - 1 \right]_+ + \sum_{r\in R} \left[ \frac{(\mathbf{w}_r^\top \mathbf{d}_r)^2}{\|\mathbf{d}_r\|_2^2} - \epsilon^2 \right]_+ \right\}$$

**TransH: 6.5M+ DrugBank Triples, $\alpha$ = 0.05, Dimensions = 100, $\lambda$ = 1, C = 0.25,**

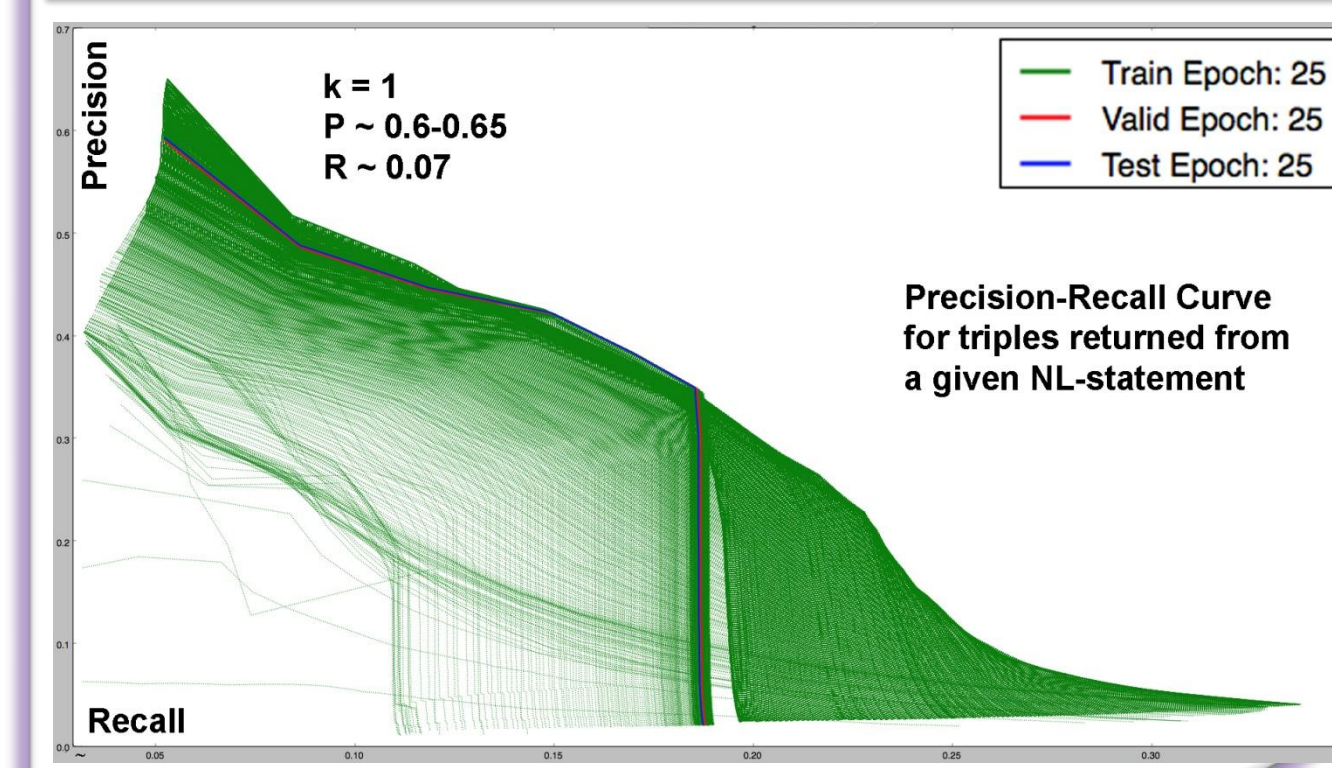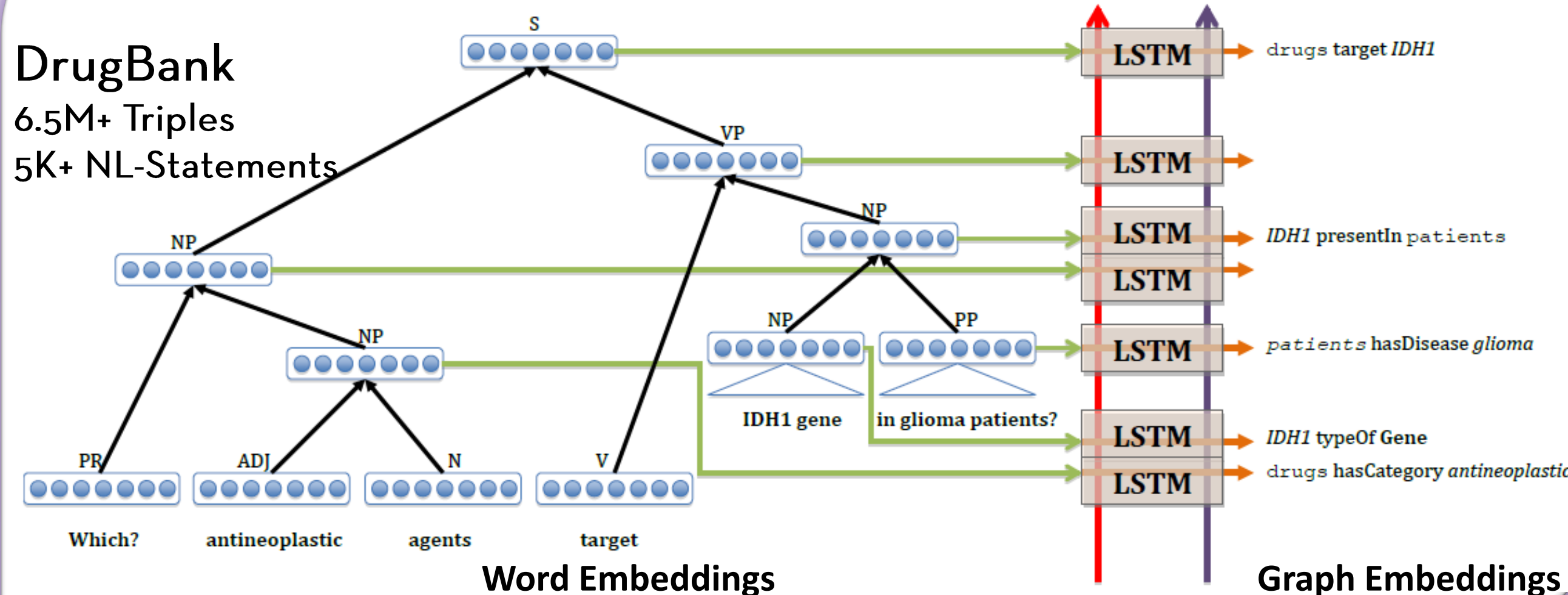**Reference:** Wang, Zhen, et al. "Knowledge Graph Embedding by Translating on Hyperplanes." *AAAI*. 2014.

## Neural Network Architecture



**DrugBank** 6.5M+ Triples 5K+ NL-Statements

drugs target *IDH1*
*IDH1* presentIn patients
patients hasDisease glioma
*IDH1* typeOf Gene
drugs hasCategory antineoplastic

IDH1 gene | in glioma patients?

Which? | antineoplastic | agents | target

**Word Embeddings**               **Graph Embeddings**

## Evaluation



**Intrinsic Evaluation**

Legend: Gene, Chemical, Disease, Single Nucleotide Polymorphism, Species, Other Mutation

**Extrinsic Evaluation: MESH Recommendations**
- Each biomedical abstract is provided with Medical Subject Headings (MESH) manually.
- **Train:** 8.3M+, **Valid:** 2.7M+, **Test:** 2.7M+ abstracts
- Baseline is KNN algorithm and simple TF-IDF vectors (K=6):- Precision: 0.31, Recall: 0.45
- 2-layer Neural Network, with 20 and 300 hidden nodes respectively, 25 Epochs, 27K MESH Terms ($\alpha$= 0.001, $\lambda$= 0.001) :- Precision: 0.86, Recall: 0.22.



k = 1
P ~ 0.6-0.65
R ~ 0.07

Train Epoch: 25
Valid Epoch: 25
Test Epoch: 25

**Precision-Recall Curve for triples returned from a given NL-statement**

## Conclusion

- Word & Graph embeddings and Neural Network architecture that can translate NL-queries to structured triples, and provide a usable interface to tackle integrative bioinformatics challenges.
- More rigorous evaluations of graph embeddings and method, and scale it to LSLOD network.