# HPOLabeler: Improving Prediction of Human Protein-Phenotype Associations by Learning to Rank

## Lizhi Liu *(Presenter)*
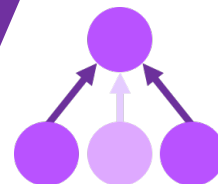
**Xiaodi Huang, Hiroshi Mamitsuka, Shanfeng Zhu**

**Email:** *liulizhi1996@gmail.com*

*School of Computer Science*
*Shanghai Key Lab of Intelligent Information Processing*
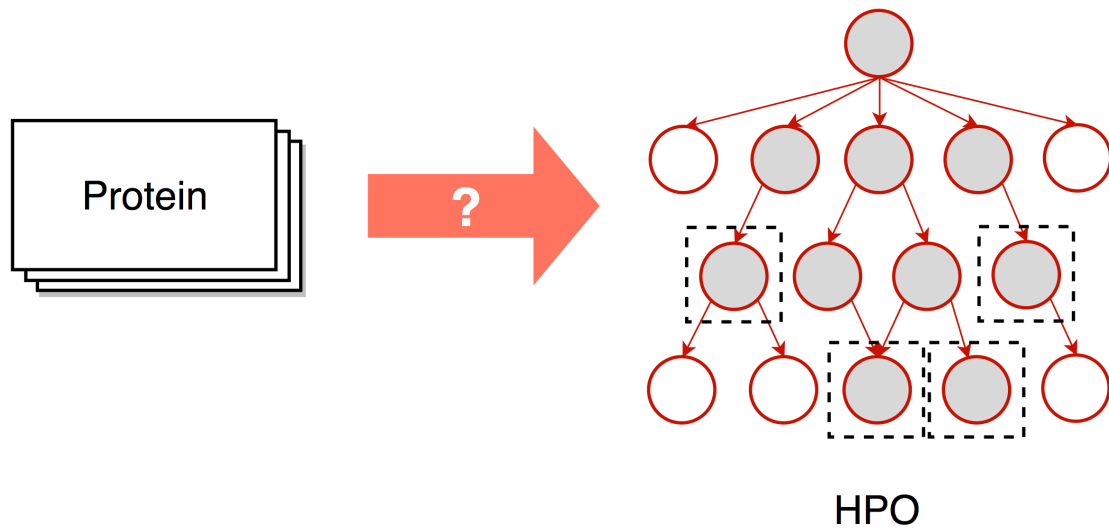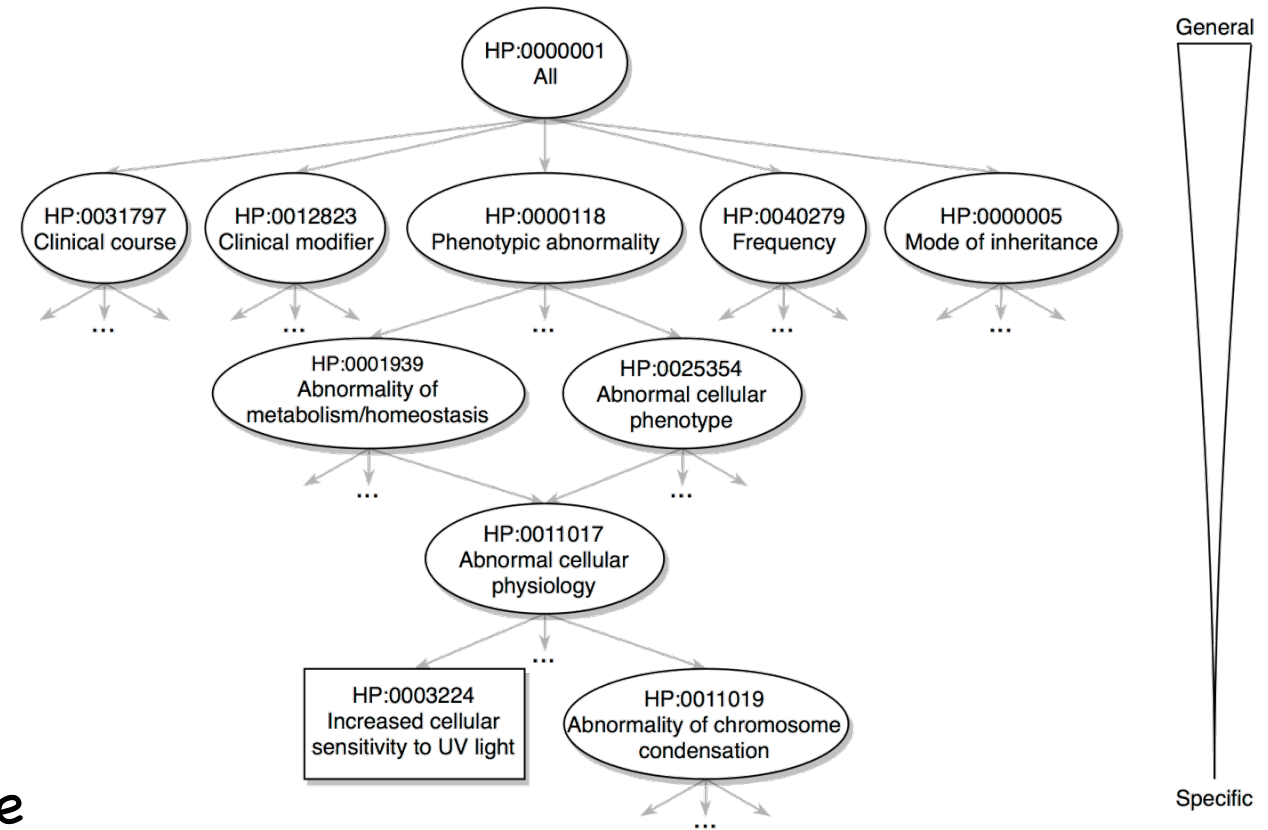*Fudan University, Shanghai, China*

**Website:** http://issubmission.sjtu.edu.cn/hpolabeler/
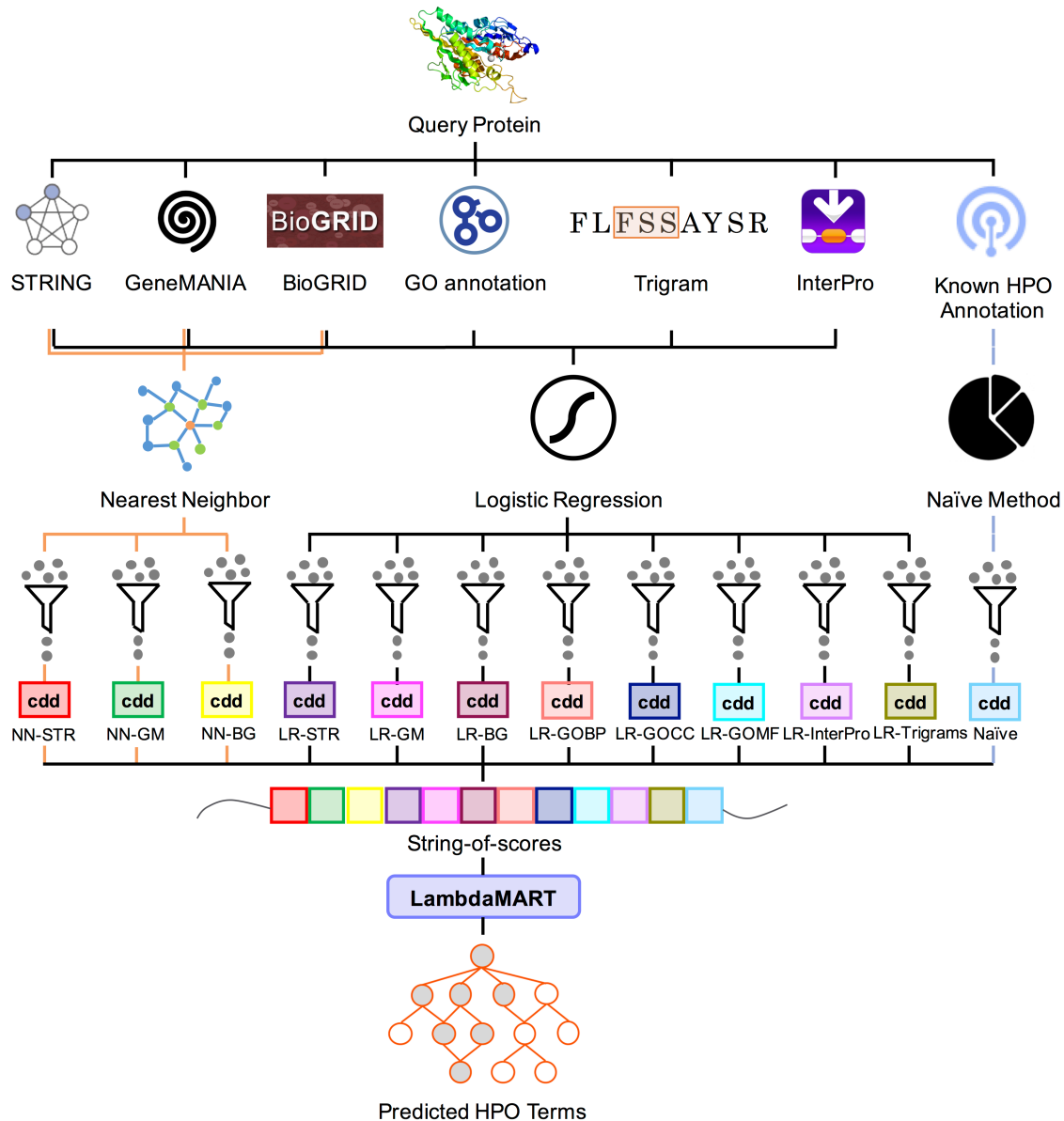
# Problem Statement



HPO

**Predicting HPO annotations of Human Proteins**

**Our goal:** Using machine learning techniques to integrate multiple data sources and improve the performance
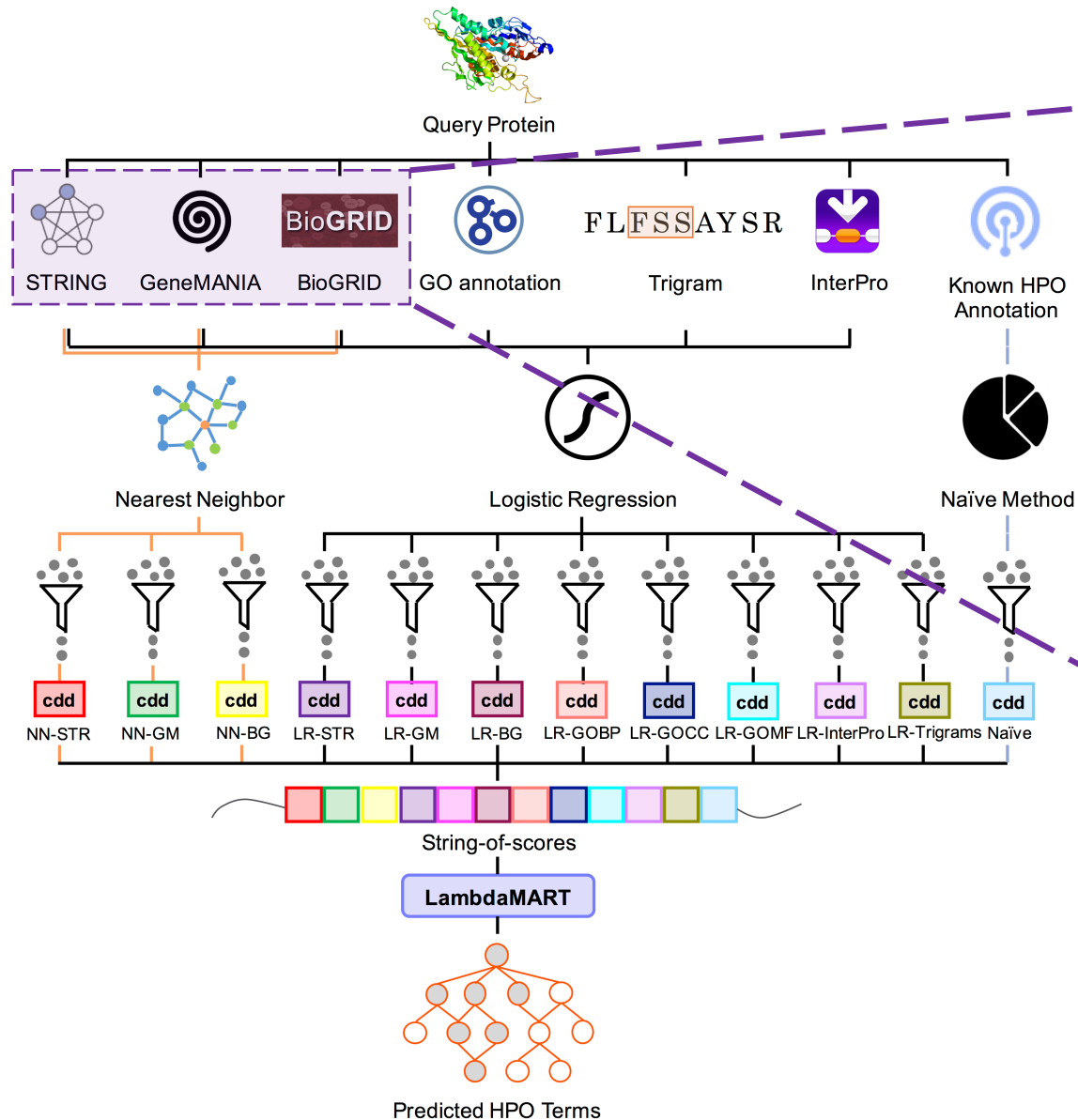
**HPO (Human Phenotype Ontology)**

# Our Proposal — HPOLabeler



## Key Points

- **Ensemble learning** : Stacking

- **Learning to Rank** to integrate multiple basic models to further improve the performance

- **Only one** better than Naïve method in temporal validation

# Feature Extraction – PPI Networks

$$\mathbf{x}_i^{(\text{STR})} = \left( x_{i,1}^{(\text{STR})}, x_{i,2}^{(\text{STR})}, \cdots, x_{i,n^{(\text{STR})}}^{(\text{STR})} \right)^T \qquad (1)$$
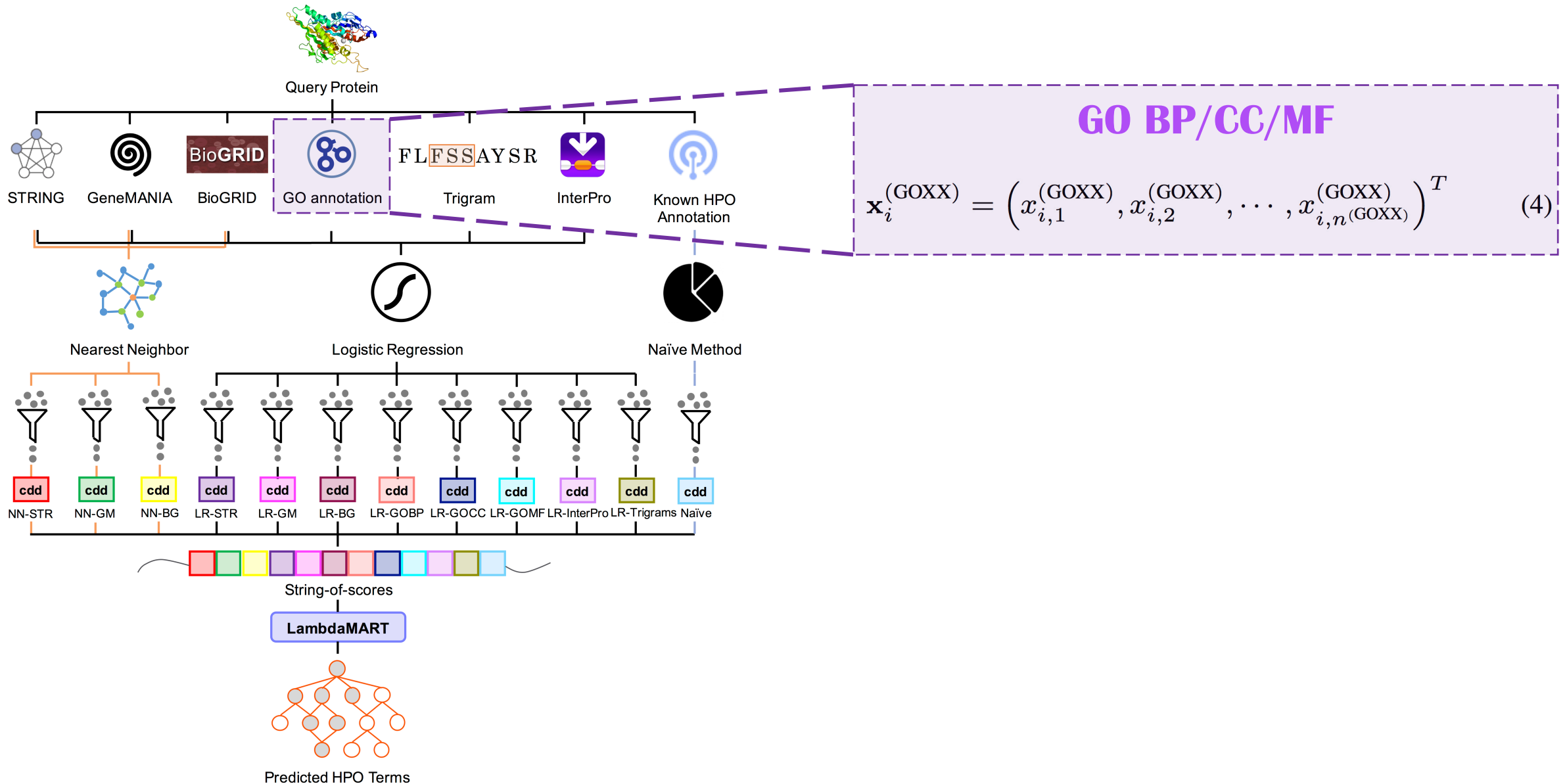
$$\mathbf{x}_i^{(\text{GM})} = \left( x_{i,1}^{(\text{GM})}, x_{i,2}^{(\text{GM})}, \cdots, x_{i,n^{(\text{GM})}}^{(\text{GM})} \right)^T \qquad (2)$$
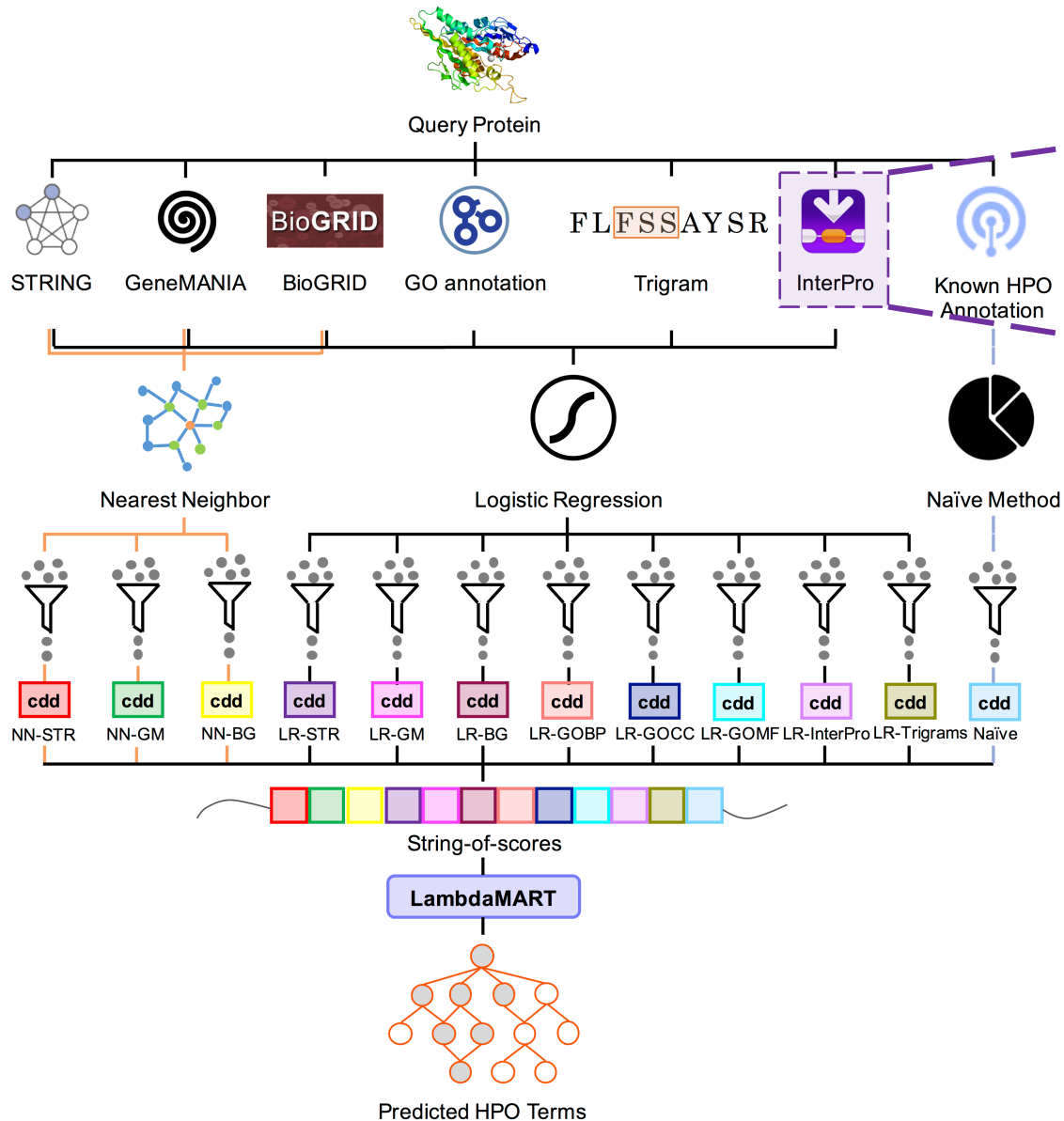
$$\mathbf{x}_i^{(\text{BGD})} = \left( x_{i,1}^{(\text{BGD})}, x_{i,2}^{(\text{BGD})}, \cdots, x_{i,n^{(\text{BGD})}}^{(\text{BGD})} \right)^T \qquad (3)$$

# Feature Extraction – GO annotations

Query Protein

STRING  GeneMANIA  BioGRID  GO annotation  Trigram  InterPro  Known HPO Annotation

Nearest Neighbor  Logistic Regression  Naïve Method

NN-STR  NN-GM  NN-BG  LR-STR  LR-GM  LR-BG  LR-GOBP  LR-GOCC  LR-GOMF  LR-InterPro  LR-Trigrams  Naïve

String-of-scores

**LambdaMART**

Predicted HPO Terms

## GO BP/CC/MF

$$\mathbf{x}_i^{(\text{GOXX})} = \left( x_{i,1}^{(\text{GOXX})}, x_{i,2}^{(\text{GOXX})}, \cdots, x_{i,n^{(\text{GOXX})}}^{(\text{GOXX})} \right)^T \qquad (4)$$
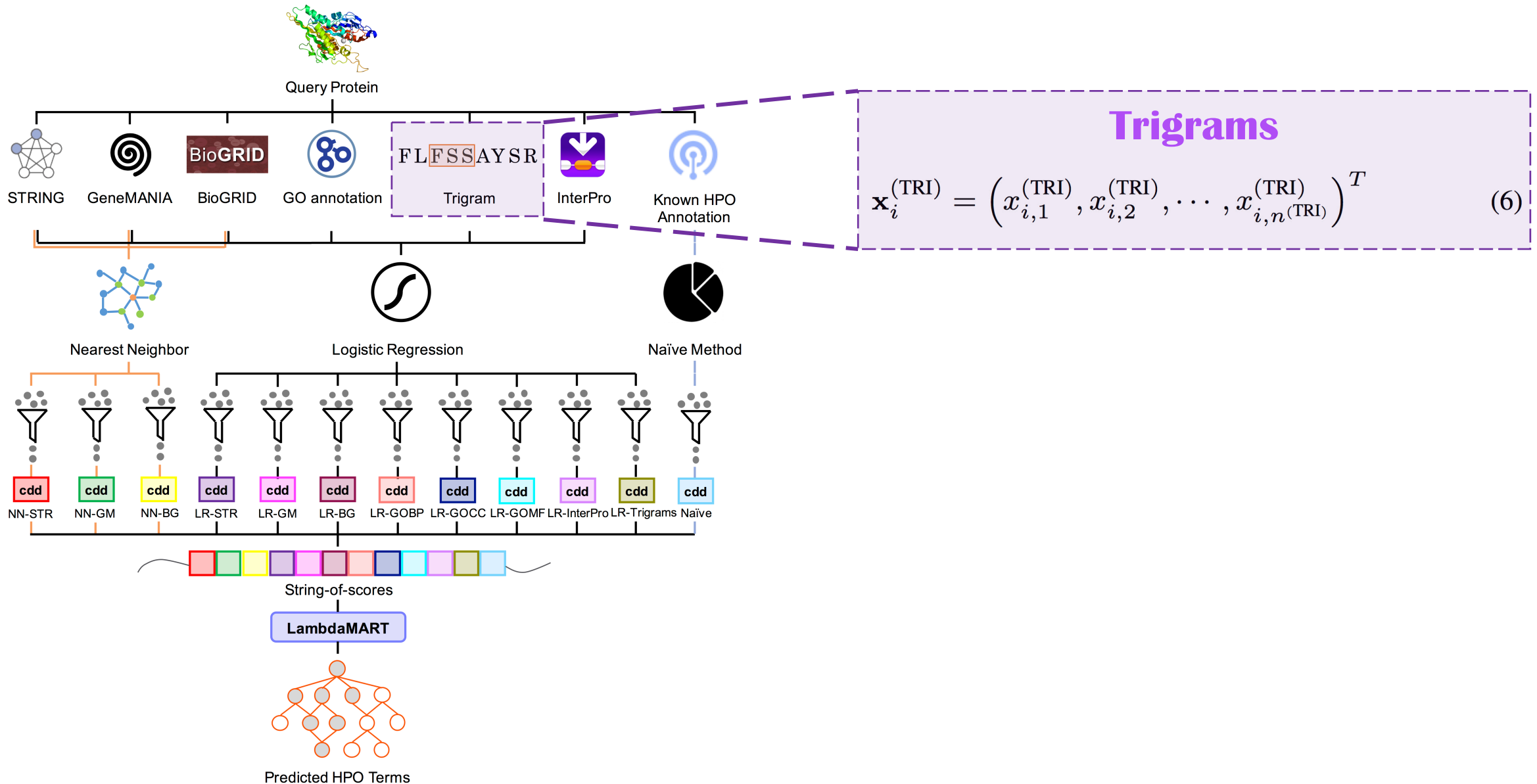
# Feature Extraction – InterPro



$$\mathbf{x}_i^{(\text{IPR})} = \left(x_{i,1}^{(\text{IPR})}, x_{i,2}^{(\text{IPR})}, \cdots, x_{i,n^{(\text{IPR})}}^{(\text{IPR})}\right)^T \tag{5}$$

InterPro signatures

# Feature Extraction – Amino Acid Sequences

Trigrams

$$\mathbf{x}_i^{(\text{TRI})} = \left( x_{i,1}^{(\text{TRI})}, x_{i,2}^{(\text{TRI})}, \cdots, x_{i,n^{(\text{TRI})}}^{(\text{TRI})} \right)^T \qquad (6)$$

# Basic Model – Logistic Regression



LR model for each HPO term

$$S^{(f)}(p_i, t) = \mathcal{L}_t^{(f)}(\mathbf{x}_i^{(f)}) = P\left(y_{i,t} = 1 | \mathbf{x}_i^{(f)}\right) \qquad (7)$$

# Basic Model – Nearest Neighbor



Query Protein

STRING  GeneMANIA  BioGRID  GO annotation  Trigram  InterPro  Known HPO Annotation

Nearest Neighbor  Logistic Regression  Naïve Method
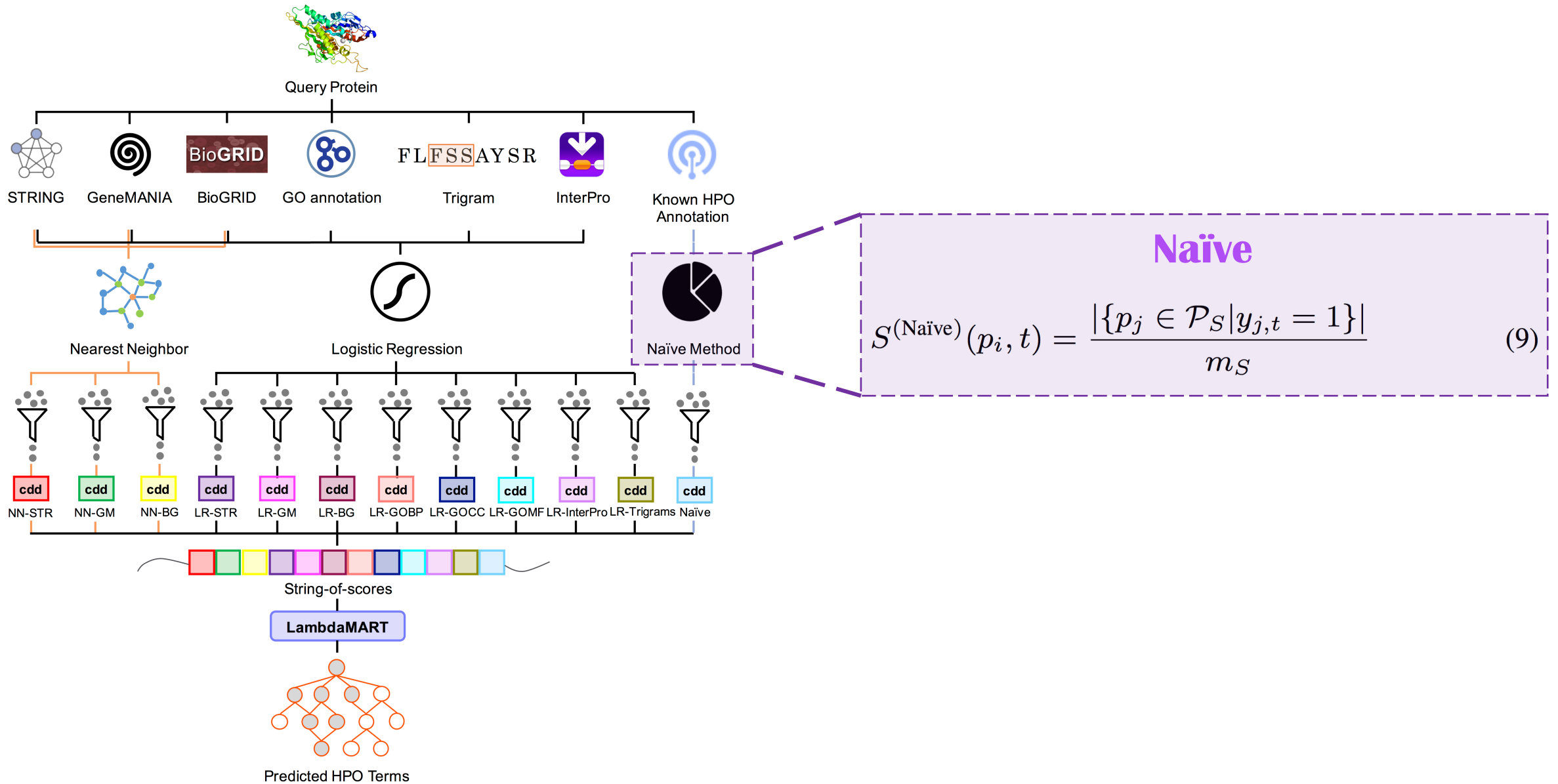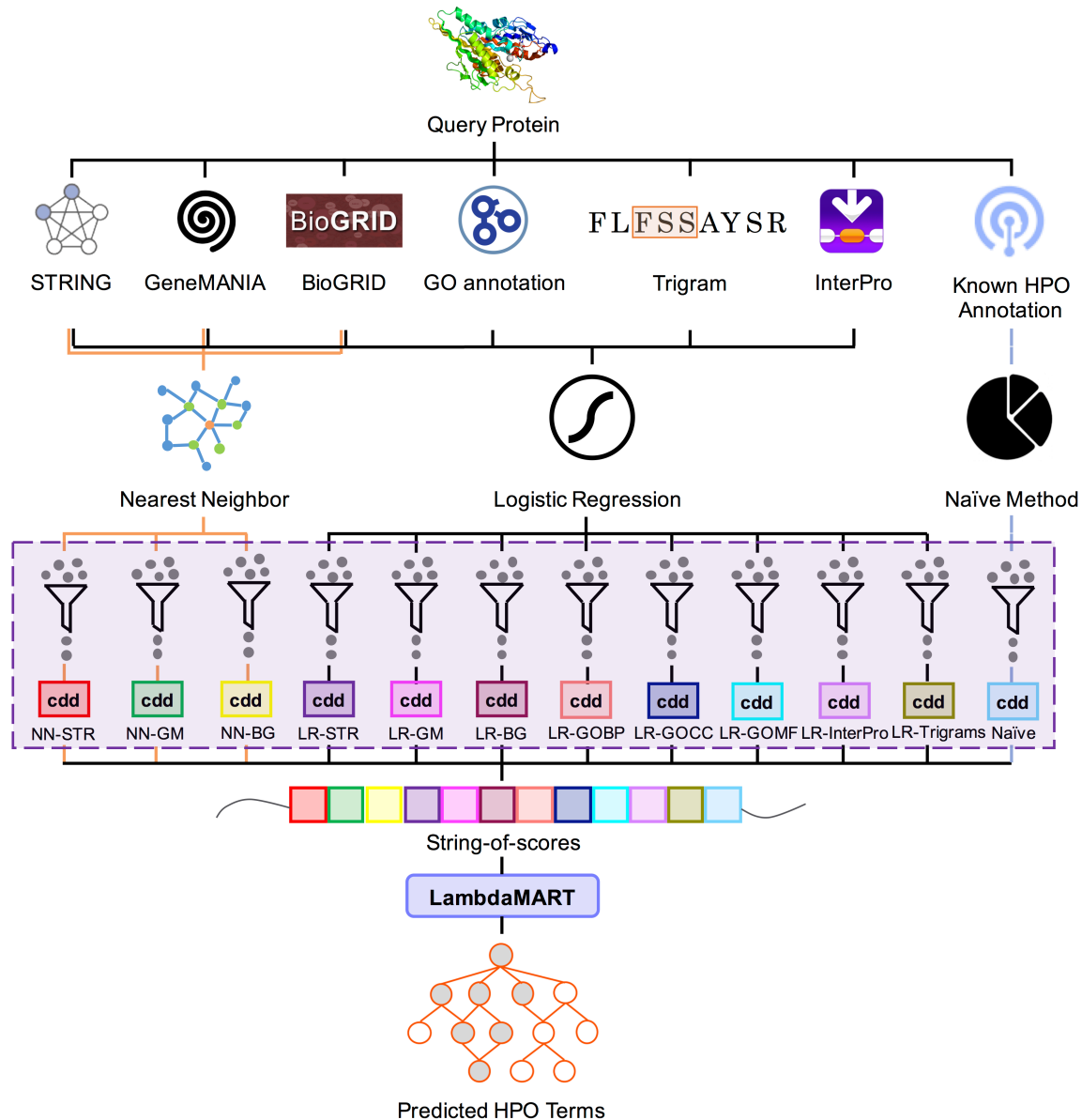
**Nearest Neighbor on STRING, GeneMANIA and BioGRID**

$$S^{(\text{NBR-G})}(p_i, t) = \frac{\sum_{p_j \in N_G(p_i)} d(p_i, p_j) \cdot y_{j,t}}{\sum_{p_j \in N_G(p_i)} d(p_i, p_j)} \qquad (8)$$
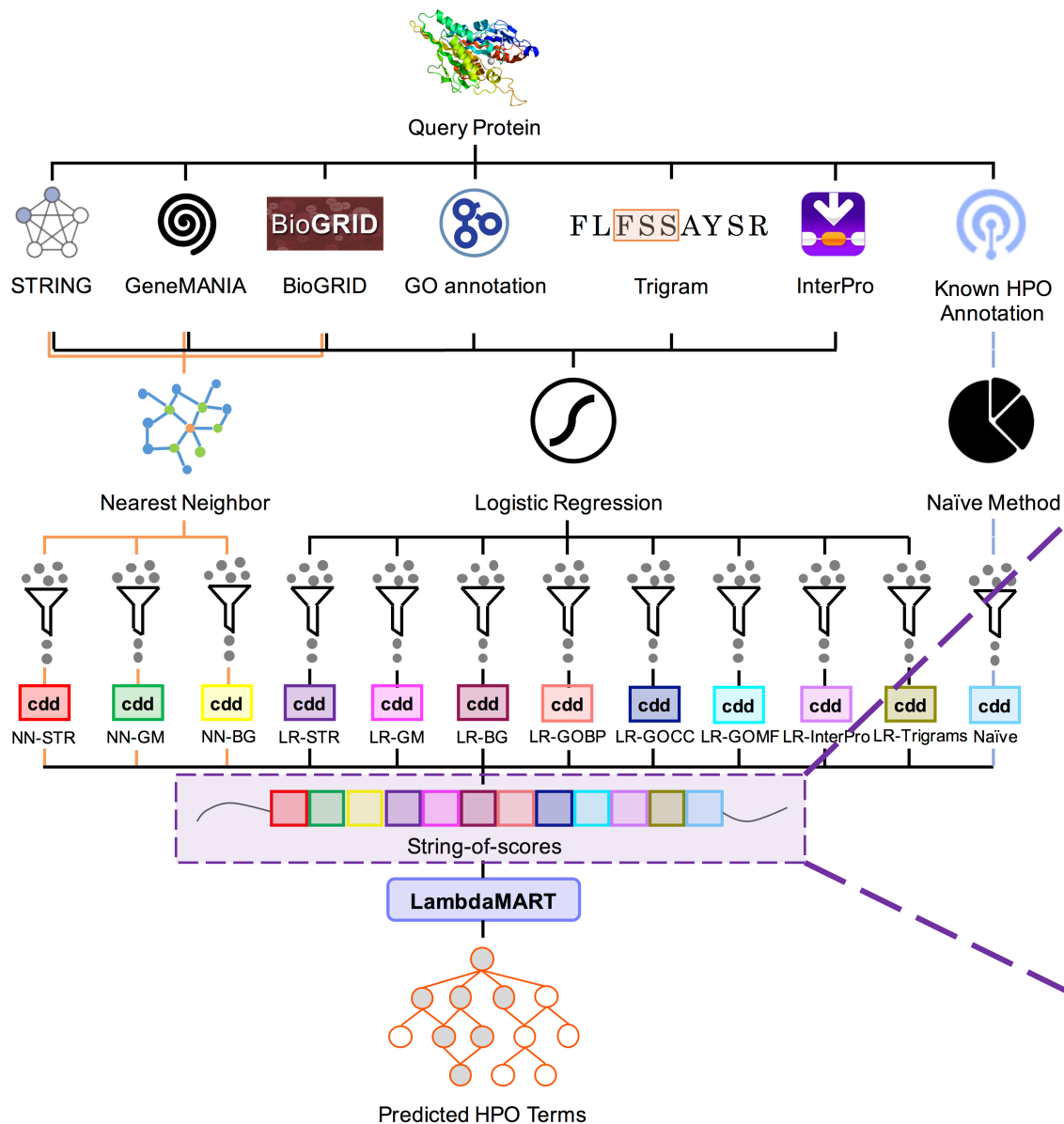
cdd  cdd  cdd  cdd  cdd  cdd  cdd  cdd  cdd  cdd  cdd  cdd

NN-STR  NN-GM  NN-BG  LR-STR  LR-GM  LR-BG  LR-GOBP  LR-GOCC  LR-GOMF  LR-InterPro  LR-Trigrams  Naïve

String-of-scores

**LambdaMART**

Predicted HPO Terms

# Basic Model – Naïve



$$S^{(\text{Naïve})}(p_i, t) = \frac{|\{p_j \in \mathcal{P}_S | y_{j,t} = 1\}|}{m_S} \qquad (9)$$

- Top-$k$ of HPO terms on each of basic models are selected
- Take the union of these subsets as the finalized candidates

# HPOLabeler – Step 3: Ranking
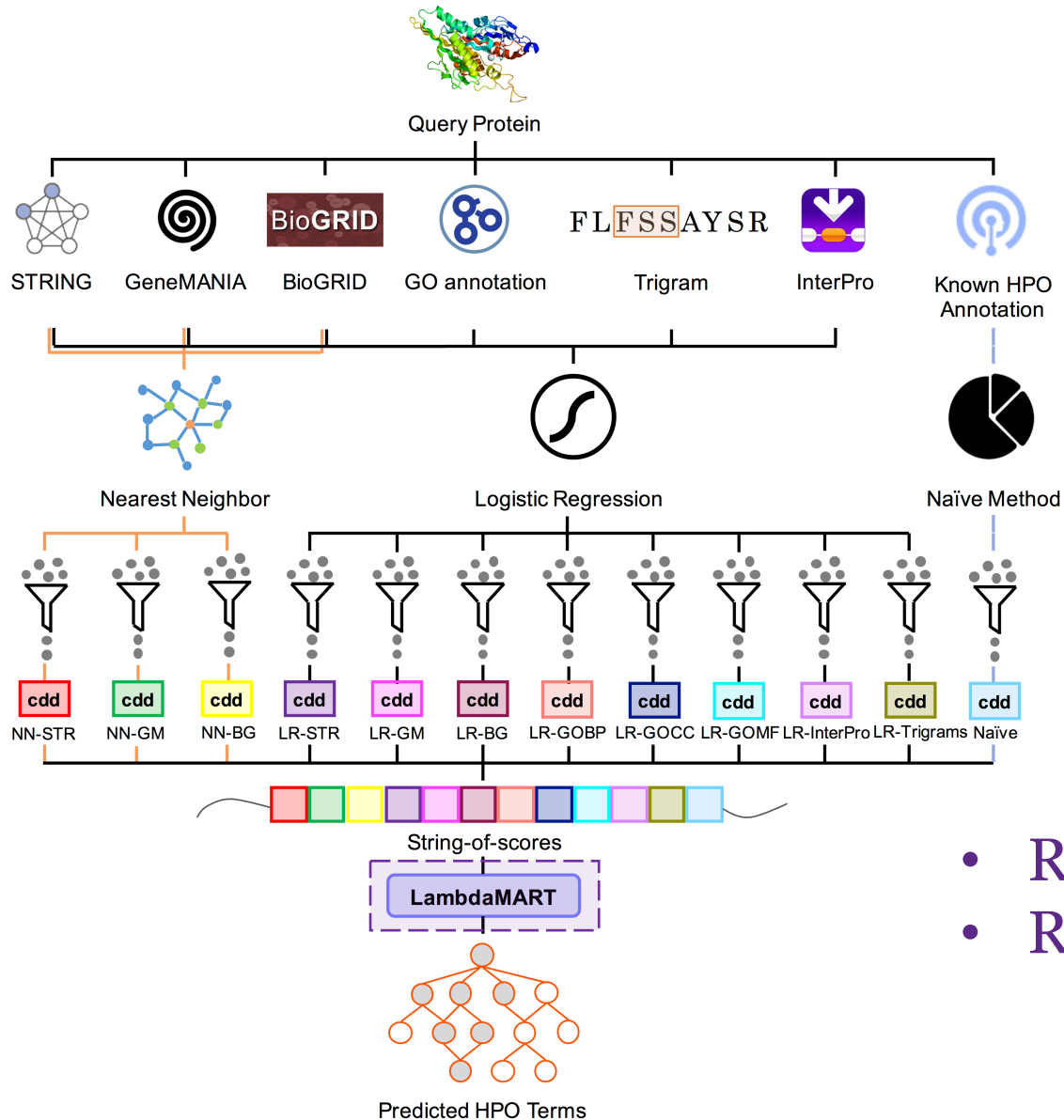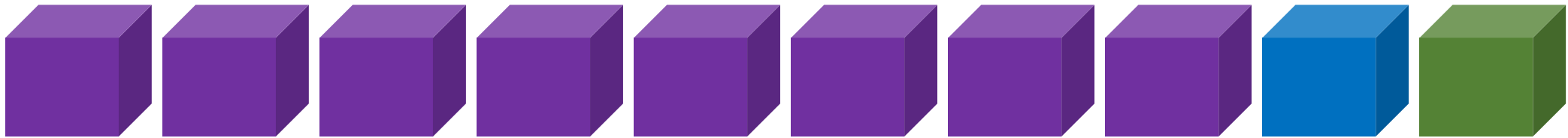


- Re-rank candidates based on **LambdaMART**
- Receive a ranked list of predictive scores

# Evaluation 1: Cross-validation

**2018-07-27**



**3,722 proteins**          **8,067 HPO terms**          **Avg. 119.4 annotations**
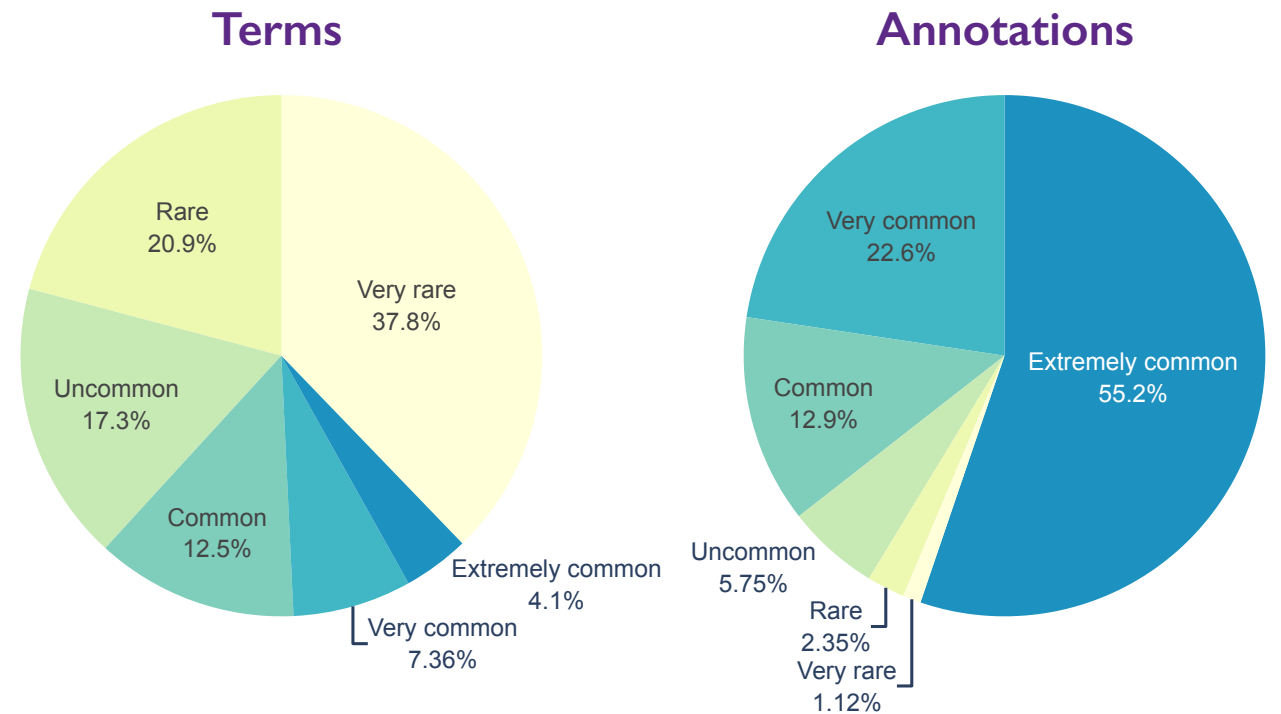
# Results of Cross-validation – Basic Models

| Component | $F_{\max}$ | AUC | AUPR |
|---|---|---|---|
| LR-STRING | 0.4174 | 0.6390 | 0.2697 |
| LR-GeneMANIA | 0.3506 | 0.7282 | 0.2605 |
| LR-BioGRID | 0.3441 | 0.5941 | 0.2677 |
| LR-GO BP | 0.3777 | 0.6741 | 0.2926 |
| LR-GO CC | 0.3643 | 0.6544 | 0.2916 |
| LR-GO MF | 0.3343 | 0.6081 | 0.2403 |
| LR-InterPro | 0.3588 | 0.6041 | 0.2699 |
| LR-Trigrams | 0.2941 | 0.5136 | 0.1564 |
| NN-STRING | **0.4213** | **0.7892** | **0.3635** |
| NN-GeneMANIA | 0.4110 | 0.7274 | 0.3550 |
| NN-BioGRID | 0.3529 | 0.6407 | 0.2822 |
| Naive | 0.3517 | 0.5 | 0.2590 |

- Nearest Neighbor 👍
- PPI 👍
- NN > LR

# Results of Cross-validation – Comparison

| Method | $F_{\max}$ | AUC | AUPR |
|---|---|---|---|
| PHENOstruct | 0.4228 | 0.7760 | 0.3596 |
| S→D→H | 0.3476 | 0.7606 | 0.2580 |
| SVM | 0.4055 | 0.6831 | 0.2900 |
| LR | 0.4242 | 0.6690 | 0.2972 |
| HTD-DAG | 0.4134 | 0.6832 | 0.2951 |
| TPR-DAG | 0.4253 | 0.6840 | 0.3170 |
| PhenoPPIOrth | 0.1430 | 0.5731 | 0.0558 |
| HPO2GO | 0.2751 | 0.5395 | 0.0936 |
| Naive | 0.3517 | 0.5 | 0.2591 |
| HPOLabeler | **0.4688**[*] | **0.7956** | **0.4293**[*] |

# Facts: HPO and Annotations are unbalance

# Results of CV – Avg. AUC group by frequency

| Method | Uncommon | Com. | Very Com. | Extremely Com. |
|---|---|---|---|---|
| PHENOstruct | **0.8161** | 0.7888 | 0.7748 | 0.7501 |
| S→D→H | 0.7925 | 0.7619 | 0.7324 | 0.6895 |
| SVM | 0.6690 | 0.6851 | 0.6989 | 0.6937 |
| LR | 0.6429 | 0.6704 | 0.6974 | 0.7023 |
| HTD-DAG | 0.6716 | 0.6842 | 0.6971 | 0.6928 |
| TPR-DAG | 0.6689 | 0.6849 | 0.7005 | 0.7009 |
| PhenoPPIOrth | 0.5961 | 0.5745 | 0.5562 | 0.5231 |
| HPO2GO | 0.5521 | 0.5347 | 0.5267 | 0.5306 |
| Naive | 0.5 | 0.5 | 0.5 | 0.5 |
| HPOLabeler | 0.7922 | **0.8046**$^{*}$ | **0.8082**$^{*}$ | **0.7778**$^{*}$ |

- High-frequency groups 😊
- Low-frequency groups 😐

# Results of CV – Leave-one-source-out



- PPI: most informative
- NN: best performing
- All changes < 0: indispensable

# Evaluation 2: Temporal Validation



|  | Train | L2R | Test |
|---|---|---|---|
| Proteins | 3,334 | 304 | 226 |
| Used HPO terms | 7,394 | 2,836 | 2,091 |
| Annotations | 107.0936 | 83.9079 | 61.5177 |

HPOLabeler Basic models Training — 2017-02-24

HPOLabeler L2R Training — 2018-03-09

HPOLabeler Test — 2018-12-21

| Method | $F_{\max}$ | AUC | AUPR |
|---|---|---|---|
| PHENOstruct | 0.3054 | 0.6362 | 0.1424 |
| S→D→H | 0.1461 | 0.5473 | 0.0603 |
| SVM | 0.2791 | 0.5929 | 0.1077 |
| LR | 0.2956 | 0.5950 | 0.1119 |
| HTD-DAG | 0.2933 | 0.5956 | 0.1138 |
| TPR-DAG | 0.3002 | 0.5962 | 0.1235 |
| PhenoPPIOrth | 0.0678 | 0.5219 | 0.0121 |
| HPO2GO | 0.2075 | 0.5083 | 0.0277 |
| Naive | 0.3097 | 0.5 | 0.2147 |
| HPOLabeler | **0.3415** | **0.6398** | **0.2342** |

# Findings: HPO annotations are incomplete

**A**



**#HPO terms associated with
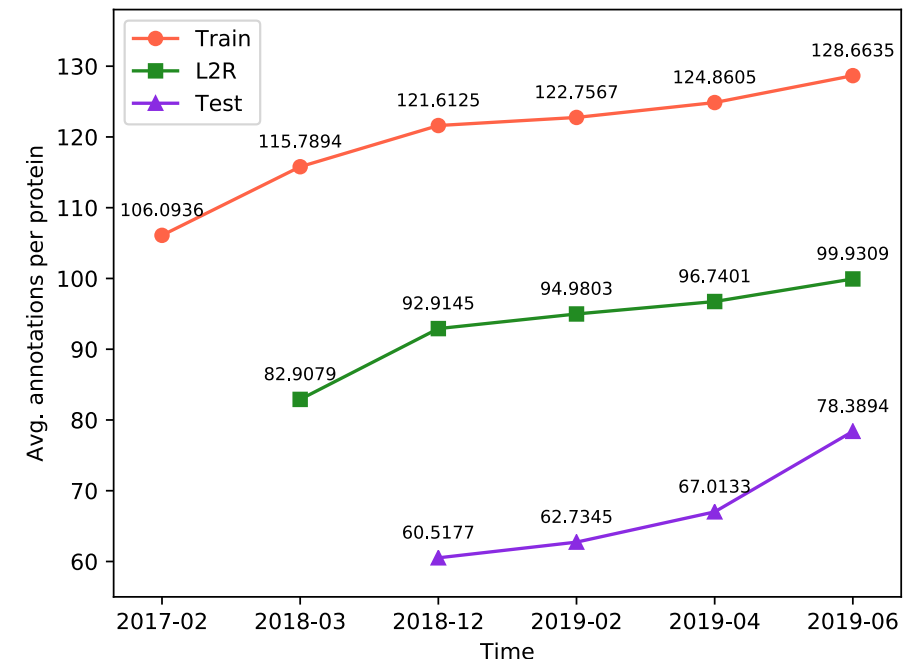a protein in each dataset**

**B**



**AUPRs evaluated by HPO annotations
released at different times**
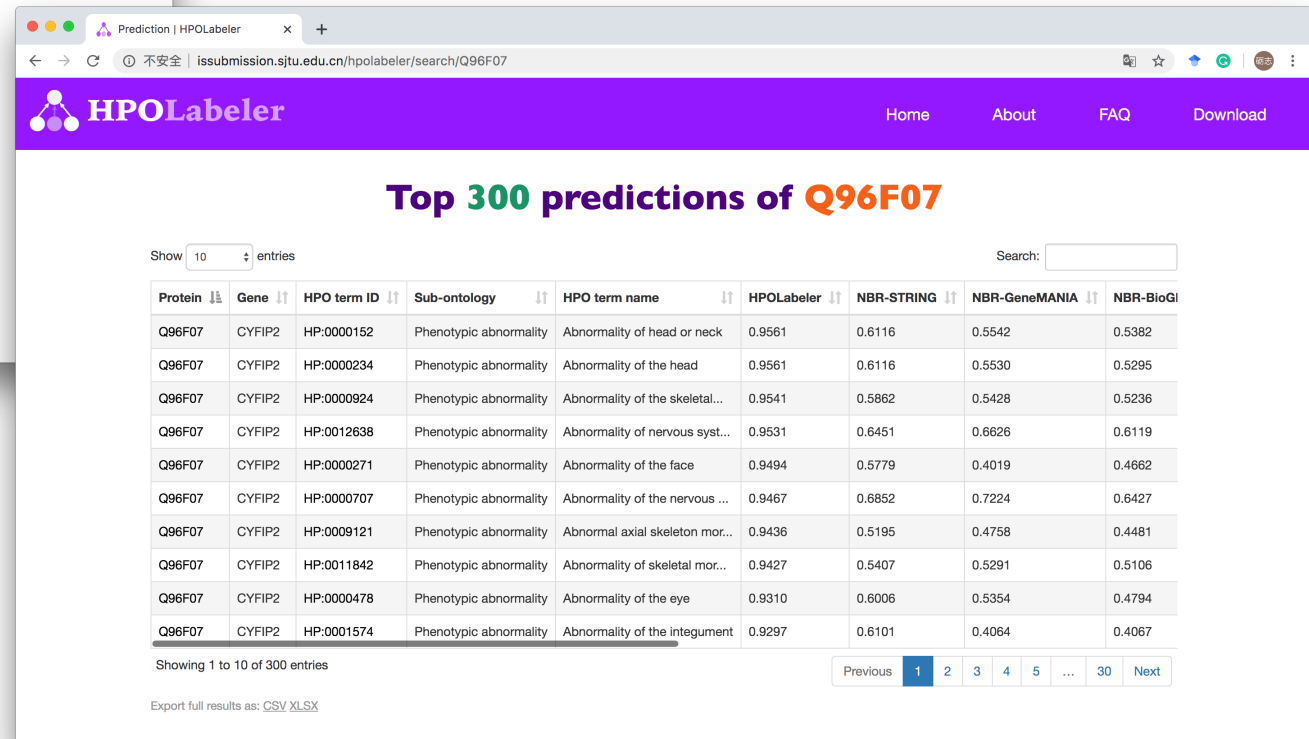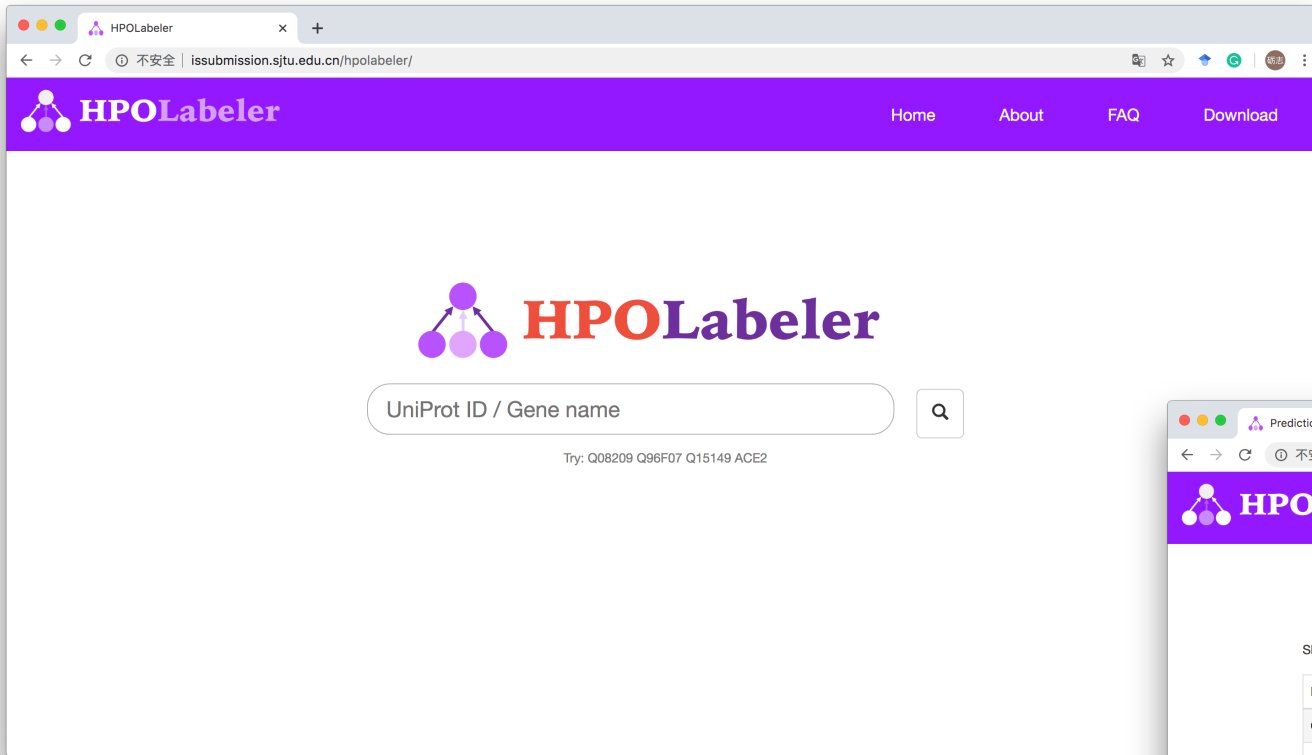
# Findings: HPO annotations are incomplete

| UniProt ID | Protein name | Gene symbol | Disease ID | HPO term ID | HPO term name | Rank |
|---|---|---|---|---|---|---|
| Q08209 | Serine/threonine-protein phosphatase 2B catalytic subunit alpha isoform | PPP3CA | ORPHA:442835 OMIM:617711 | HP:0000924 | Abnormality of the skeletal system | 3 |
| | | | | HP:0011842 | Abnormality of skeletal morphology | 9 |
| | | | | HP:0025031 | Abnormality of the digestive system | 18 |
| Q96F07 | Cytoplasmic FMR1-interacting protein 2 | CYFIP2 | ORPHA:442835 OMIM:618008 | HP:0000152 | Abnormality of head or neck | 1 |
| | | | | HP:0000234 | Abnormality of the head | 1 |
| | | | | HP:0000924 | Abnormality of the skeletal system | 3 |
| P61981 | 14-3-3 protein gamma | YWHAG | ORPHA:442835 OMIM:617665 | HP:0000478 | Abnormality of the eye | 3 |
| | | | | HP:0000152 | Abnormality of head or neck | 8 |
| | | | | HP:0000234 | Abnormality of the head | 9 |

**Predicted associations *(Excerpt)* which were evaluated as negatives by old annotations but appeared in the latest release in Feb. 2019, meaning that all are actually positives**

**Avg. #HPO annotations of newly added proteins keep increasing with time**

# Online Platform



http://issubmission.sjtu.edu.cn/hpolabeler/

# Conclusions

- We propose HPOLabeler, which is able to integrate diverse types of evidences including PPI, GO, InterPro and trigrams, in the framework of Learning to Rank.

- We empirically validated the performance of HPOLabeler, which significantly outperformed all competing methods.

- Further examinations of the experimental results indicate that:
  - PPI is the most informative data source;
  - lower predictive performance in temporal validation might be caused by incomplete annotations of new proteins.

- We developed an online platform:
  http://issubmission.sjtu.edu.cn/hpolabeler/

# THANK YOU