# Building Gene Networks by Information Extraction, Cleansing, & Integration

**Limsoon Wong**

NUS
National University
of Singapore

---

NUS
National University
of Singapore

## Plan

- **Motivation for study of gene network**
- **Example Efforts at I2R**
  - Disease Pathweaver
  - Dragon ERG Solution
- **Technical Challenges Involved**
  - Name entity recognition
  - Co-reference resolution
  - Protein interaction extraction
- **Discussion on issues**

Motivation

---

# Some Useful Information Sources

- **Disease-centric resources**
  - OMIM [NCBI OMIM, 2004]
  - Gene2Disease [Perez-Iratxeta et al., *Nature*, 2002]
  - MedGene [Hu et al.. *J. Proteome Res.*, 2003]
- **Emphasized direct gene-disease relationships**
  - Provide lists of disease-related genes
  - Do not provide info on gene-gene interactions & their networks

- **Related interaction resources**
  - KEGG [Kanehisa, *NAR*, 2000]
- **Manually constructed protein interaction networks**
- **Mostly metabolic pathways and few disease pathways**
  - Only 7 disease pathways

Adapted w/ permission from Zhou Zhang, Suisheng Tang, & See-Kiong Ng

# Why Gene Network?

**NUS**
National University
of Singapore

- **Many common diseases are**
  - not caused by a genetic variation within a single gene
- **But are influenced by**
  - complex interactions among multiple genes
  - environmental & lifestyle factors

**Monogenic → Heterogenic**

Adapted w/ permission from Zhuo Zhang, Suisheng Tang, & See-Kiong Ng

---

# Desired Outcome of Gene Network Study

**NUS**
National University
of Singapore

- **Help scientists understand the mechanism of complex diseases by**
  - Greatly reducing work load for primary study of genetic diseases, broaden the scope of molecular studies
  - Easily identifying key players in the gene network, help in finding potential drug targets
- **Scalability framework**
  - Extend to many genetic diseases
  - Include other resources of gene interactions

Adapted w/ permission from Zhou Zhang, Suisheng Tang, & See-Kiong Ng

Some Gene Network Study Efforts at I²R

**Disease Pathweaver**
**Dragon ERG Solution**

---

# Disease Pathweaver, Zhang et al. *APBC 2005*

- **Automatic constructing disease pathways**
  - Identify core genes
  - Mine info on core genes
  - Construct interaction network betw core genes
- **Data sources:**
  - Online literature
  - High-thru'put biological data



Adapted w/ permission from Zhou Zhang, Suisheng Tang, & See-Kiong Ng

Disease Pathweaver:
# The Portal

- **Portal for human nervous diseases gene networks**
  - http://research.i2r.a-star.edu/NSDPath

- **Statistics**
  - 37 Human Nervous System Disorders
  - 7 ~ 60 core genes per disease
  - 2 ~ 320 core interactions per disease

Adapted w/ permission from Zhou Zhang, Suisheng Tang, & See-Kiong Ng

Copyright © 2006 by Limsoon Wong. Adapted w/ permission from Zhuo Zhang, Suisheng Tang, & See-Kiong Ng

Copyright © 2006 by Limsoon Wong. Adapted w/ permission from Zhuo Zhang, Suisheng Tang, & See-Kiong Ng



Copyright © 2006 by Limsoon Wong. Adapted w/ permission from Zhuo Zhang, Suisheng Tang, & See-Kiong Ng

Disease Pathweaver: A Tour

Copyright © 2006 by Limsoon Wong. Adapted w/ permission from Zhuo Zhang, Suisheng Tang, & See-Kiong Ng



Disease Pathweaver:
DPW vs KEGG on Huntington Disease

Adapted w/ permission from Zhuo Zhang, Suisheng Tang, & See-Kiong Ng

# Estrogen-Responsive Genes

**NUS** National University of Singapore

- **Why**
  - Affects human physiology in many aspects
  - Related to many diseases
  - Widely used in clinic
- **Challenges**
  - Multiple pathways
  - Difficult to predict ERE
  - Many estrogen-responsive genes but only a few are well-studied
  - Difficult to keep up w/ speed of knowledge accumulation

- **Needed**
  - Tools to predict ERE & estrogen-responsive genes
  - Database of useful info
  - Systems to predict impt regulatory units, associate gene functions, & generate global view of gene network

Adapted w/ permission from Suisheng Tang & Vlad Bajic

---

# Dragon ERG Solution

**NUS** National University of Singapore

**ERE dependent**

→ E2 dependent

→ E2 independent

ERE Finder
ERG Finder

ERGDB

ERG Explorer

**ERE independent**

→ ERs bind to other TFs

→ Membrane receptors

Coming soon!

*text mining here!*

Adapted w/ permission from Suisheng Tang & Vlad Bajic

## Slide 1

Dragon ERG Solution:
# Dragon ERE Finder,
Bajic et al, *NAR*, 2003



➢Predict functional ERE in genomic DNA
➢One prediction in 13.3k bp
➢Allow further analysis

```
Sequence name:          gi|18590884:438190-434690|chromosome|19
Sequence length:        3501
# of bases:             A=906, C=974, G=928, T=693
Expected ERE sensitivity: 83%

Predicted positions relative to the 5' end of the input sequence
ERE (red)

Forward strand
1363 ○ AG-GGTCA-CTT-CGGCC-CA new pattern
3301 ◉ AA-GATCA-GTC-TGGCC-AA new pattern
# of ERE guesses = 2

Reverse complement strand
# of ERE guesses = 0
```

**http://sdmc.i2r.a-star.edu.sg/ERE-V2/index**

## Slide 2

Dragon ERG Solution:
# Dragon Estrogen-Responsive Gene Finder, Tang et al, *NAR*, 2004a



➢Only for human genome
➢Using 117 bp ERE frame
➢Evaluated by PubMed & microarray data

**http://sdmc.i2r.a-star.edu.sg/DRAGON/ERGP1_0**

Dragon ERG Solution:

# Dragon Estrogen-Responsive Genes Database, Tang et al, *NAR,* 2004b

Dragon ERG Solution

BCL2

➢Contains >1000 genes
➢Manually curated
➢Basic gene info
➢Experimental evidence
➢Full set of references
➢ERE sites annotated

**Gene Information**

| Organism: | Homo sapiens (human) |
|---|---|
| Gene: | BCL2 |
| Alternate Symbols: | Bcl-2 |
| Description: | B-cell CLL/lymphoma 2 |
| Chromosome: | 18 |
| Contig: | NT_025028 |
| Locus Link: | 596 |
| Unigene Entry: | Hs.79241 |
| Refseq/GenBank Link: | M13995, NM_000633, M13994, M14745, BC027258, X06487, NM_000657 |

**Experimental Information**

| Method: | RT-PCR, Northern, Southern and Western blots |
|---|---|
| In vivo/In vitro: | in vitro |
| Regulation: | up |
| Cell Line/Tissue: | ovary, surface epithelium cell lines |
| Time Point: | 24 h - 6 days |
| | PUID:11356682 |

**http://sdmc.i2r.a-star.edu.sg/promoter/Ergdb-v11**

Dragon ERG Solution:
## DEERGF

Dragon ERG Solution

**Dragon Explorer of Estrogen Responsive Genes Functionality (DEERGF)**
Version 1.0

**ATF3 (Homo sapiens)**

⦿ Link to ERGDB ver 1.1 for gene **ATF3** (Homo sapiens)
○ Orthologs of gene **ATF3** (Homo sapiens)
○ *Ab-initio* DNA motifs found in the ortholog group of gene **ATF3** (Homo sapiens)
○ TFs and GO categories associate via text-mining to gene **ATF3**
○ Gene expression (using eVOC) of **ATF3** (Homo sapiens)

Submit    Reset

**Graphical Representation for Ortholog Gene**

| ATF3 (Hs) | |
| Atf3 (Mm) | |
| Atf3 (Rn) | |

Dragon ERG Solution:
## Case-Specific TF relation networks, Pan et al, *NAR*, 2004

- **Analyse abstracts**
- **Stemming, POS tagging**
- **Use ANNs, SVM, discriminant analysis**
- **Simplified rules for sentence analysis**
- **Constraints on the forms of sentences**
- **Sensitivity ~75%**
- **Precision ~82%**

- **Produce reports & direct links to PubMed docs, & graphical presentations of entity links**

Adapted w/ permission from Vlad Bajic

---

## Technical Challenges

**Named entity recognition**

**Co-reference resolution**

**Data cleansing**

**NUS**
**National University of Singapore**

# Bio Entity Name Recognition,

### Zhou et al., *BioCreAtIvE,* 2004

**LEGEND**

- Virus
- Tissue
- RNA
- Protein
- Polynucleotide
- Peptide
- OtherOrganicCompound
- OtherName
- OtherArtificialSource
- Organism
- Nucleotide
- MultiCell
- MonoCell
- Lipid
- Inorganic
- DNA
- CellType
- CellLine
- CellComponent
- Carbohydrate
- BodyPart
- Atom
- AminoAcidMonomer

Erythropoietin stimulates transcription of the TAL1/SCL gene and phosphorylation of its protein products.

Activation of the TAL1 (or SCL) gene, originally identified through its involvement by a recurrent chromosomal translocation, is the most frequent molecular lesion recognized in T-cell acute lymphoblastic leukemia. The protein products of this gene contain the basic-helix-loop-helix motif characteristic of a large family of transcription factors that bind to the canonical DNA sequence CANNTG as protein heterodimers. TAL1 expression by erythroid cells in vivo and in chemical-induced erythroleukemia cell lines in vivo suggested the gene might regulate aspects of erythroid differentiation. Since the terminal events of erythropoiesis are controlled by the glycoprotein hormone erythropoietin (Epo), we investigated whether the expression or activity of the TAL1 gene and its protein products were affected by Epo in splenic erythroblasts from mice infected with an anemia-inducing strain of Friend virus (FVA cells). Epo elicited a rapid, dose-related increase in TAL1 mRNA by increasing transcription of the gene and stabilizing one of its mRNAs. An Epo-inducible TAL1 DNA binding activity was identified in FVA cell nuclear extracts that subsequently decayed despite accumulating mRNA and protein. Induction of DNA binding activity was associated temporally with Epo-induced phosphorylation of nuclear TAL1 protein. These results indicate that Epo acts at both transcriptional and posttranscriptional levels on the TAL1 locus in Friend virus-induced erythroblasts and establish a link between Epo signaling mechanisms and a member of a family of transcription factors involved in the differentiation of diverse cell lineages.

Adapted w/ permission from Jian Su

---

### Bio Entity Name Recognition:
# Ensemble Classification Approach

- **Features considered**
  - orthographic, POS, morphologic, surface word, trigger words ($TW_1$: receptor, enhancer, etc. $TW_2$: activation, stimulation, etc.)
- **SVM**
  - Context of 7 words
  - Each word gives 5 features, plus its position
- **HMMs**
  - 3 features used (orthographic, POS, surface word)
  - $HMM_1$ & $HMM_2$ use POS taggers trained on diff corpora

$HMM_1$
balanced precision & recall

$HMM_2$
low precision & high recall

SVM
high precision & low recall

Majority Voting

12

Bio Entity Name Recognition:

# Performance at *BioCreAtIvE* 2004

| Modules | Closed-1 | Closed-2 | Closed-3 | Open-1 |
|---|---|---|---|---|
| SVM | Surface word, orthographic feature, suffix, trigger | | | |
| | GENIA-POS | **Refined-BioCreative-POS** | Refined-BioCreative-POS | Refined-BioCreative-POS |
| HMM1 | Surface word, orthographic feature, | | | |
| | GENIA-POS | **Refined-BioCreative-POS** | Refined-BioCreative-POS | Refined-BioCreative-POS |
| HMM2 | Surface word, orthographic feature, BioCreative-POS | | | |
| Ensemble | Majority Voting | | | |
| Abbreviation Res. | Abbreviation Resolution based on the parentheses structure | | | |
| Refinement of protein/gene names | N/A | N/A | **YES** | N/A |
| Dictionary Matching | Closed Dictionary | Closed Dictionary | Closed Dictionary | **Open Dictionary** |
| Overall Performance | P79.97 R80.15 F80.23 | P80.46 R80.80 F80.63(**+0.40**) | P82.00 R83.17 F82.58(**+2.35**) | P75.10 R81.26 F78.06(**-4.52**) |

Adapted w/ permission from Zhou et al., BioCreAtIve 2004

---

# Co-Reference Resolution,

Yang et al., *IJCNLP*, 2004

During training, for each anaphor NPj in a given text, a positive instance is generated by pairing NPj with its closest non-pronominal antecedent. A set of negative instances is also formed by NPj and each of the non-pronominal markables occurring between NPj and NPi.

A training instance is associated with a feature vector which, as described in Table 2, consists of 16 features, 12 of which are used in Soon *et al.*'s system. Here two string match features are tried in

When the training instances are ready, a classifier is learned by C5.0 algorithm (Quinlan, 1993).

During resolution, each encountered noun phrase, NPj, is paired in turn with each preceding noun phrase, NPi, from right to left. Each pair is associated with a feature vector as during training, and then presented to the coreference classifier. The classifier returns a positive or negative result indicating whether or not NPi is coreferential to NPj. The process terminates once an antecedent is found for NPj, or the beginning of the text is reached. In the former case, NPj is to be linked into the coreferential chain where the antecedent occurs

Adapted w/ permission from Yang et al., IJCNLP 2004

Co-Reference Resolution:

# Baseline Features Used

**NUS** National University of Singapore

**Features describing the antecedent candidate (NPi):**

| 1. | ante_DefNp | 1 if NPi is a definite NP; else 0 |
| 2. | ante_DemoNP | 1 if NPj start with a demonstrative; else 0 |
| 3. | ante_IndefNP | 1 if NPi is an indefinite NP; else 0 |
| 4. | ante_Pron | 1 if NPi is a pronoun; else 0 |
| 5. | ante_ProperNP | 1 if NPi is a proper NP; else 0 |

**Features describing the anaphor candidate (NPj):**

| 6. | ana_DefNP | 1 if NPj is a definite NP; else 0 |
| 7. | ana_DemoNP | 1 if NPj start with a demonstrative; else 0 |
| 8. | ana_IndefNP | 1 if NPj is an indefinite NP; else 0 |
| 9. | ana_Pron | 1 if NPj is a pronoun; else 0 |
| 10. | ana_ProperNP | 1 if NPj is a proper NP; else 0 |

**Features describing the antecedent candidate (NPi) and the possible anaphor (NPj):**

| 11. | GenderAgree | 1 if NPi and NPj agree in gender; else 0 if disagree; -1 if unknown |
| 12. | NumAgree | 1 if NPi and NPj agree in number; 0 if disagree; -1 if unknown |
| 13. | Appositive | 1 if NPi and NPj are in an appositive structure; else 0 |
| 14. | Alias | 1 if NPi and NPj are in an alias of the other; else 0 |
| 15 | SemanticAgree | 1 if NPi and NPj agree in semantic class; 0 if disagree; -1 if unknow |
| | | 1 if NPi and NPj contain the same head string; else 0 |
| 16. | HeadStrMatch | 1 if NPi and NPj contain the same string after discarding determiners; |
| 16' | FullStrMatch | else 0 |

Table 2: Feature set for the baseline coreference resolution system

Adapted w/ permission from Yang et al., IJCNLP 2004

---

Co-Reference Resolution:

# New Features Used & Performance

**NUS** National University of Singapore

**Features describing the antecedent candidate (NPi):**

| 17. | ante_Relative | 1 if NPi is modified by a relative clause; else 0 |
| 18. | ante_specialNP | 1 if NPi is a special definite np which acts as a non-anaphor; else 0 |

**Features describing the anaphor candidate (NPj):**

| 19. | ana_Relative | 1 if NPj modified by a relative clause; else 0 |
| 20. | ana_SpecialNP | 1 if NPj is a special definite np which refers to no antecedent; else 0 |

**Features describing the antecedent candidate (NPi) and the possible anaphor (NPj):**

| 21- | ante_ana_(EntireNP, Number, Verb, Prep, | Matching degree of NPi.(EntireNP, ..., Com- |
| 29. | AdjJ, AdjR, AdjS, ProperNP, CommonNP) | monNP) and NPj.(EntireNP, ..., CommonNP) |
| | | |
| 30- | ana_ante_(EntireNP, Number, Verb, Prep, | Matching degree of NPj.(EntireNP, ..., Com- |
| 38 | AdjJ, AdjR, AdjS, ProperNP, CommonNP) | monNP) and NPi.(EntireNP, ..., CommonNP) |

Table 4: New string matching features of our coreference resolution system

| | Recall | Precision | F-measure |
|---|---|---|---|
| *HeadStrMatch* | **71.4** | 53.1 | 60.9 |
| *FullStrMatch* | 51.0 | 68.5 | 58.4 |
| *NewFeature\** | 70.5 | 63.8 | 66.9 |
| *NonAnaphor+NewFeature\** | 68.1 | **69.7** | **68.9** |

**Base Classifier: C5.0**

Table 5: Experimental results on the Medline data set using C5.0 (the *ed systems use *ContainRa-tion* metric with *Binary* weighting scheme)

Adapted w/ permission from Yang et al., IJCNLP 2004

# Protein Interaction Extraction,

Xiao et al, IJCNLP 2004

Tokenization and Morphological Analysis

↓

POS Tagging

↓

Name Entity Recognition

↓

Sentence Analysis

↓

Co-reference Resolution

↓

Maximum Entropy Classifier

- **The max entropy model:**

$$P(o, h) = \frac{1}{Z(h)} \prod_{j=1}^{k} \alpha_j^{f_j(h,o)}$$

- **where**
  - o is the outcome
  - h is the feature vector
  - Z(h) is normalization function
  - $f_j$ are feature functions
  - $\alpha_j$ are feature weights

Adapted w/ permission from Xiao et al., IJCNLP 2004

---

Protein Interaction Extraction:

# Features Used

| Feature names | Feature values |
|---|---|
| First protein name | p1_bovine, p1_prion, p1_protein |
| Second protein name | p2_protein, p2_kinase |
| Words between two protein names | b_strongly, b_interact, b_with, b_the, . |
| Left words | l_here, l_that, l_recombine |
| Right words | r_. |
| Overlap | ProteinNameInBetween=0 |
| Keyword | Keyword=interacts_between |
| Chunk heads in between | chunk_head_strongly, chunk_head_interacts, chunk_head_with, chunk_head_alpha/alpha', chunk_head_subunit, chunk_head_of |
| Surrounding chunk heads | leftChunkHead=here_that, rightChunkHead=interacts |
| Chunk types in between | ChunkType=ADVP_VP_PP_NP_NP_PP |
| Parser tree path | PaserPath=NPB_S_VP_PP_NP_PP |
| Dependent | Dependent=false |
| Dependent root | DependentRoot=interacts, DependentRootPos=VBZ |
| Pair of two protein heads | PairOfProteinHead=prion_kinase |
| Pair of abbreviations | AbbreviationPair=bprp_protein_kinase |

Adapted w/ permission from Xiao et al.

Protein Interaction Extraction:
# Performance on IEPA Corpus

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Words in two names | * | * | * | * | * | * | * | * | * | * |
| Words between two names | * | * | * | * | * | * | * | * | * | * |
| Surrounding words | | * | * | * | * | * | * | * | * | * |
| Overlap | | * | | | | | | | | |
| Keyword feature | | | | * | * | * | * | * | * | * |
| Chunk features | | | | | * | * | * | * | * | * |
| Parse tree | | | | | | * | * | * | * | |
| Dependent tree | | | | | | | * | * | * | |
| Pair of proteins | | | | | | | | * | * | * |
| Abbreviation pair | | | | | | | | | * | * |
| Recall (%) | 80.5 | 86.1 | 85.9 | 86.6 | 87.2 | 87.1 | 87.2 | 90.1 | 93.6 | 93.9 |
| Precision (%) | 75.0 | 81.2 | 81.1 | 81.7 | 83.1 | 83.0 | 82.8 | 85.3 | 88.0 | 88.0 |
| F-measure | 77.5 | 83.6 | 83.3 | 84.1 | 85.1 | 85.0 | 84.9 | 87.7 | 90.7 | 90.9 |

IEPA = Interaction Extraction Performance Assessment

Adapted w/ permission from Xiao et al.

---

# Some Interesting Issues in Constructing Gene/Protein Networks

## Issues

- **Sound:**
  - Is the contents of our databases correct?
- **Complete:**
  - Is the structure of our databases expressive enough to capture critical information explicitly?
- **Understandable:**
  - Is our databases or search results understandable?
- **Other issues relating to NLP/IE**

---

Soundness:
Is the content of our databases correct?



**NUS**
National University
of Singapore

# Sources of Errors, Koh et al, *DBiDB*, 2005

- **11 types & 28 subtypes of data artifacts**
  - Critical artifacts (vector contaminated sequences, duplicates, sequence structure violations)
  - Non-critical artifacts (misspellings, synonyms)
- **> 20,000 seq records in public contain artifacts**
- **Identification of these artifacts are impt for accurate knowledge discovery**

- **Sources of artifacts**
  - Diverse sources of data
    - **Repeated submissions of seqs to db's**
    - **Cross-updating of db's**
  - Data Annotation
    - **Db's have diff ways for data annotation**
    - **Data entry errors can be introduced**
    - **Different interpretations**
  - Lack of standardized nomenclature
    - **Variations in naming**
    - **Synonyms, homonyms, & abbrevn**
  - Inadequacy of data quality control mechanisms
    - **Systematic approaches to data cleaning are lacking**

Adapted w/ permission from Judice Koh, Mong Li Lee, & Vladimir Brusic

# Classification of Errors

Adapted w/permision from Judice Koh, Mong Li Lee, & Vladimir Brusic

18

Completeness:
Is the structure of our
databases expressive enough
to capture critical information
explicitly?

**NUS**
National University
of Singapore

---

**NUS**

## Expressive Power

- **Take a key paper such as the Kohn paper that summarises current knowledge on p53 regulation.**
- **Is there a structured database that is able to capture all info in that paper explicitly?**
- **Is there a semi-structured database that is able to capture all info in that paper explicitly?**
- **How well does this (semi-) structured database generalize to other similar type of papers?**

Understandability:
Is our databases or
search results
understandable?

---

# Self-Organization

- **Take a search on p53. You will get >300k hits or some number like that on MEDLINE**
- **It is not feasible for anyone to go thru all of that to find what he wants!** **And this problem is growing bigger as MEDLINE doubles every 1-2 year.**
- **Need to organize the database and/or the search results into hierarchy or "semantic" net to make it easier for users to understand or to browse the results**
- **How do we define this hierarchy/net?**
- **Can this hierarchy/net be self-organized?**

# Problems relating to NLP/IE

## Handling full-length papers

- **Source document structure parsing**
- **Hyper-linked file tracking**
- **Figure and table processing**
- **Special symbol handling**

# Information retrieval

- **Document and sentence retrieval**
- **Relevant interaction filtering**

# Bio name recognition

- **Nomenclature loosely followed**
- **Frequent use of conjunction and disjunction in bio names with multiple bio-entity names sharing one head noun**
- **Long descriptive names**
- **Names of genes and proteins used interchangeably**

# Bio-interaction extraction

- **Inherent complexity of biological interactions**

  Generally, biological interactions involve both first-order basic molecular interaction events:

  `<molecule_1> <interact> <molecule_2>`

  as well as complex second-order causal events such as:

  `<event_1>` *is_caused_by* `<event_2>`
  *provided* `<event_3>`

- ⇒ **Sentences describing them also tend to be complicated**

---

# Bio-interaction extraction

- **Domain knowledge is often needed for interaction template filling**



"c-Abl tyrosine kinase activity is blocked by pRb, which binds to the c-Abl kinase domain."

(pRb *inhibit* tyrosine kinase *activity-of* c-Abl)
(pRb *bind-to* c-Abl kinase domain)

(pRb *inhibit* tyrosine kinase *activity-of* c-Abl)
*is-caused-by*
(pRb *bind-to* c-Abl *at* kinase domain)

# Extraction of other relevant info

- **Contextual information**
  - Species, cell type, cellular localisation, etc
- **Negative information**

> negation sentence such as "We have found no evidence that protein A is involved in the regulation of gene B." is often reported

- **Speculative & incomplete facts**

> "We suggest that HNF-3 may play a dual role on glucagon gene transcription by 1) inhibiting the transactivation potential of Pax6 on the G1 and G3 elements and 2) direct activation through G1 and G2."

---

# Information integration

- **Bio-name mapping**

> "IL-1, alone, or in combination with IFN-gamma and TNF-alpha leads to islet cell dysfunction and death."
> (4)

> "Exposure of flourescence activated cell sorting (FACS)-purified rat and mouse beta cells to interleukin-1beta (IL-1beta), in combination with IFN-gamma and/or TNF-alpha, leads to cell death by necrosis and predominantly by apoptosis."
> (5)

- **Bio-interaction mapping**
  - how do you know two complex sentences are talking about the same interaction?

## Resource for training & benchmarking

- **Is there such a good resource, especially for the more complex tasks?**

---

## Acknowledgements



Data Cleansing:
Judice Koh, Vladimir Brusic,
Mong Li Lee, Asif M. Khan,
Paul T.J. Tan, Heiny Tan,
Kenneth Lee, Wilson Goh,
Songsak Tongchusak,
Kavitha Gopalakrishnan

Pathweaver:
Zhuo Zhang, See-Kiong Ng,
Suisheng Tang, Chris Tan

Dragon ERG & TF Miner:
Suisheng Tang, Vlad Bajic,
Zuo Li, Pan Hong,
Vidhu Chaudhary, Raja Kangasa

Info Extraction:
Guodong Zhou, Jian Su,
ChewLim Tan, Juan Xiao,
Xiaofeng Yang, Chris Tan,
Dan Shen, Jie Zhang

25

## References

- Zhang et al., "Toward discovering disease-specific gene networks from online literature", *APBC*, 3:161-169, 2005
- Pan et al., "Dragon TF association miner: A system for exploring transcription factor associations through text mining", *NAR*, 32:W230-W234, 2004
- Tang et al., "Computational method for discovery of estrogen-responsive genes", *NAR*, 32:6212-6217, 2004a
- Tang et al., "ERGDB: Estrogen-responsive genes database", *NAR*, 32:D533-D536, 2004b
- Bajic et al., "Dragon ERE finded ver.2: A tool for accurate detection and analysis of estrogen-response elements in vertebrate genomes", *NAR*, 31:3605-3607, 2003
- Koh et al., "A Classification of Biological Data Artifacts", *DBiBD*, 2005

## References

- Zhou et al., "Recognition of protein and gene names from text using an ensemble of classifiers and effective abbreviation resolution", *Proc. BioCreAtIvE Workshop*, pp 26-30, 2004
- Yang et al., "Improving Noun Phrase Co-reference Resolution by Matching Strings", *IJCNLP*, 1:226-333, 2004
- Xiao et al., "Protein-protein interaction extraction: A supervised learning approach", submitted