

# Annotated Texts and Resources for Written Language

Junichi TSUJII

National Centre for  
Text Mining

Computer Science  
University of Manchester

UK

Computer Science  
University of Tokyo

JAPAN

# Speech and Written Text

- Speech
  - Observables
  - Obvious applications inside speech proper, Speech recognition, Speech synthesis
- Written Text
  - Non-observables, Semantics/Pragmatics
  - Machine Translation (inside written text)
  - Information Extraction, Semantic web, e-Publishing, Terminology, Ontology and Knowledge Management, Intelligent IR and Meta-data Extraction

- **Multi-Linguality**
- **Multi-Domains**
  - Biology/Medicine/Health
  - Financial Markets
  - Legal documents
- **Multi-Tasks**
  - Machine Translation, CLIR
  - Information Extraction, Automatic Summarization, etc.
  - Ontology and Knowledge Management, Semantic Web
  - E-Journal, Knowledge Discovery

# Text Annotation

## • Task-oriented Annotation

- Bio-Creative annotated text

- System development

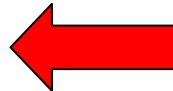
- Defined by specific tasks

- Specific curation tasks in specific environments
- Mapping of Protein names to database IDs in specific text types
- Specific event types such as Protein-Protein Interaction, in specific text types
- Disease-Gene Association of specific sets of diseases

Market failures



**Interoperable Tools**



## • Task-neutral Annotation

- GENIA Corpus [U-Tokyo, NaCTeM]

Development of generic tools

- Defined by theories

- Linguistics
  - Tokens
  - POS
  - Phrase Structure
  - Dependency Structure
  - Deep Syntax (PAS)
- Biology/Medicine/Health
  - Named Entities of various semantic types
  - Events
- Linguistics + Biology/Medicine/Health
  - Co-references

# Text Annotation

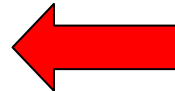
- **Task-oriented Annotation**

- Bio-Creative annotated text
- System development
- Defined by specific tasks

- Specific curation tasks in specific environments
- Mapping of Protein names to database IDs in specific text types
- Specific event types such as Protein-Protein Interaction, in specific text types
- Disease-Gene Association of specific sets of diseases



## Interoperable Tools



- **Task-neutral Annotation**

- GENIA Corpus [U-Tokyo, NaCTeM]

Development of generic tools

- Defined by theories

- **Linguistics**

- Tokens
- POS
- Phrase Structure
- Dependency Structure
- Deep Syntax (PAS)

- **Biology/Medicine/Health**

- Named Entities of various semantic types
- Events

- **Linguistics + Biology**

- Co-references

# Text Annotation

- Task-oriented Annotation

- Bio-Creative annotated text
- System development
- Defined by specific tasks

- Specific curation tasks in specific environments
- Mapping of Protein names to database IDs in specific text types
- Specific event types such as Protein-Protein Interaction, in specific text types
- Disease-Gene Association of specific sets of diseases

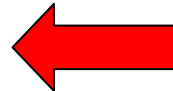


## Interoperable Tools

Development of generic tools

- Task-neutral Annotation

- GENIA Corpus [U-Tokyo, NaCTeM]



- Defined by theories

- Linguistics

- Tokens
- POS
- Phrase Structure
- Dependency Structure
- Deep Syntax (PAS)

- Biology/Medicine/Health

- Named Entities of various semantic types
- Events

- Linguistics + Biology

- Co-references

How far should we go from linguistics to conceptual contents ?

# Annotation of GENIA corpus – Term&POS

PMID:1984449

Induction of NF-KB during monocyte differentiation by HIV type 1 infection.

PMID:1984449

Induction<sub>NN</sub> of<sub>IN</sub> NF-KB<sub>NN</sub> during<sub>IN</sub> monocyte<sub>NN</sub> differentiation<sub>NN</sub> by<sub>IN</sub> HIV<sub>NN</sub> type<sub>NN</sub> 1<sub>CD</sub> infection<sub>NN</sub>.PERIOD

The<sub>DT</sub> production<sub>NN</sub> of<sub>IN</sub> human<sub>U</sub> immunodeficiency<sub>NN</sub> virus<sub>NN</sub> type<sub>NN</sub> 1<sub>CD</sub> (HIV-1<sub>NN</sub>)<sub>PRE</sub> progeny<sub>NN</sub> was<sub>VED</sub> followed<sub>VEN</sub> in<sub>IN</sub> the<sub>DT</sub> U937<sub>NN</sub> promonocytic<sub>U</sub> cell<sub>NN</sub> line<sub>NN</sub> after<sub>IN</sub> stimulation<sub>NN</sub> either<sub>CC</sub> with<sub>IN</sub> retinoic<sub>U</sub> acid<sub>NN</sub> or<sub>CC</sub> PMAN<sub>NN</sub>.COMMA and<sub>CC</sub> in<sub>IN</sub> purified<sub>VEN</sub> human<sub>U</sub> monocytes<sub>NNS</sub> and<sub>CC</sub> macrophages<sub>NNS</sub>.PERIOD Electrophoretic<sub>U</sub> mobility<sub>NN</sub> shift<sub>NN</sub> assays<sub>NNS</sub> and<sub>CC</sub> Southwestern<sub>NN</sub> blotting<sub>NN</sub> experiments<sub>NNS</sub> were<sub>VED</sub> used<sub>VEN</sub> to<sub>TC</sub> detect<sub>VE</sub> the<sub>DT</sub> binding<sub>NN</sub> of<sub>IN</sub> cellular<sub>U</sub> transactivation<sub>NN</sub> factor<sub>NN</sub> NF-KB<sub>NN</sub> to<sub>TC</sub> the<sub>DT</sub> double<sub>U</sub> repeat-KB<sub>U</sub> enhanc<sub>NN</sub> sequence<sub>NN</sub> located<sub>U</sub> in<sub>IN</sub> the<sub>DT</sub> long<sub>U</sub> terminal<sub>U</sub> repeat<sub>NN</sub>.PERIOD PMAN<sub>NN</sub> treatment<sub>NN</sub>.COMMA and<sub>CC</sub> not<sub>RE</sub> retinoic<sub>U</sub> acid<sub>NN</sub> treatment<sub>NN</sub> of<sub>IN</sub> the<sub>DT</sub> U937<sub>NN</sub> cells<sub>NNS</sub> acts<sub>VEZ</sub> in<sub>IN</sub> inducing<sub>VEG</sub> NF-KB<sub>NN</sub> expression<sub>NN</sub> in<sub>IN</sub> the<sub>DT</sub> nuclei<sub>NNS</sub>.PERIOD In<sub>IN</sub> nuclear<sub>U</sub> extracts<sub>NNS</sub> from<sub>IN</sub> monocytes<sub>NNS</sub> or<sub>CC</sub> macrophages<sub>NNS</sub>.COMMA induction<sub>NN</sub> of<sub>IN</sub> NF-KB<sub>NN</sub> occurred<sub>VED</sub> only<sub>RE</sub> if<sub>IN</sub> the<sub>DT</sub> cells<sub>NNS</sub> were<sub>VED</sub> previously<sub>RE</sub> infected<sub>VEN</sub> with<sub>IN</sub> HIV-1<sub>NN</sub>.PERIOD When<sub>WER</sub> U937<sub>NN</sub> cells<sub>NNS</sub> were<sub>VED</sub> infected<sub>VEN</sub> with<sub>IN</sub> HIV-1<sub>NN</sub>.COMMA not

Term  
and  
20  
abstracts

repeat-KB<sub>U</sub> enhanc<sub>NN</sub> sequence<sub>NN</sub> located<sub>U</sub> in<sub>IN</sub> the<sub>DT</sub> long<sub>U</sub> terminal<sub>U</sub> repeat<sub>NN</sub>.PERIOD PMAN<sub>NN</sub> treatment<sub>NN</sub>.COMMA and<sub>CC</sub> not<sub>RE</sub> retinoic<sub>U</sub> acid<sub>NN</sub> treatment<sub>NN</sub> of<sub>IN</sub> the<sub>DT</sub> U937<sub>NN</sub> cells<sub>NNS</sub> acts<sub>VEZ</sub> in<sub>IN</sub> inducing<sub>VEG</sub> NF-KB<sub>NN</sub> expression<sub>NN</sub> in<sub>IN</sub> the<sub>DT</sub> nuclei<sub>NNS</sub>.PERIOD In<sub>IN</sub> nuclear<sub>U</sub> extracts<sub>NNS</sub> from<sub>IN</sub> monocytes<sub>NNS</sub> or<sub>CC</sub> macrophages<sub>NNS</sub>.COMMA induction<sub>NN</sub> of<sub>IN</sub> NF-KB<sub>NN</sub> occurred<sub>VED</sub> only<sub>RE</sub> if<sub>IN</sub> the<sub>DT</sub> cells<sub>NNS</sub> were<sub>VED</sub> previously<sub>RE</sub> infected<sub>VEN</sub> with<sub>IN</sub> HIV-1<sub>NN</sub>.PERIOD When<sub>WER</sub> U937<sub>NN</sub> cells<sub>NNS</sub> were<sub>VED</sub> infected<sub>VEN</sub> with<sub>IN</sub> HIV-1<sub>NN</sub>.COMMA not

# Annotation of GENIA corpus –

## Process&Tree

Tree  
annotation  
2000  
abstracts

PMID:MEDLINE:1984449

NP

NP Induction

PP of

NP NF-KB

PP during

NP monocyte differentiation

PP by

NP HIV type 1 infection

PMID:1984449

Induction of NF-KB<sub>T9</sub> during monocyte<sub>T9</sub> differentiation<sub>T9</sub> by HIV type 1<sub>T9</sub> infection<sub>T9</sub>.

INFECTION P1

THEME: [T9]

PREDICATE: Induction of NF-KB during monocyte differentiation by HIV type 1 infection.

INDUCTION P2

AGENT: [I1]

THEME: [P1]

PREDICATE: Induction of NF-KB during monocyte differentiation by HIV type 1 infection.

The production of human immunodeficiency virus type 1<sub>T9</sub> (HIV-1<sub>T9</sub>) progeny<sub>T9</sub> was

Process  
annotation  
500 abstracts  
by May 2006  
1000 abstracts  
by Dec. 2006



# Task-neutral Annotation

- Tool Development
  - Training, Development, Test
  - **Domain Adaptation**

# Tool1: POS Tagger

The peri-kappa B site mediates human immunodeficiency

DT NN NN NN VBZ JJ NN

virus type 2 enhancer activation in monocytes ...

NN NN CD NN NN IN NNS

- General-Purpose POS taggers, trained by WSJ
  - Brill's tagger, TnT tagger, MX POST, etc.
  - 97%
- General-Purpose POS taggers do not work well for MEDLINE abstracts

# Errors seen in TnT tagger (Brants 2000)

A chromosomal translocation in ...

DT JJ NN IN

... and membrane potential after mitogen binding.

CC NN NN IN NN ~~JJ~~

... two factors, which bind to the same kappa B enhancers...

CD NNS WDT ~~NN~~ TO DT JJ NN NN NNS

... by analysing the Ag amino acid sequence.

IN VBG DT ~~VBG~~ JJ NN NN

... to contain more T-cell determinants than ...

TO VB ~~RBR~~ ~~JJ~~ NNS IN

Stimulation of interferon beta gene transcription in vitro by

NN IN JJ JJ NN NN ~~IN~~ ~~NN~~ IN

# Performance of GENIA Tagger

• GENIA tagger

(Ref.) TnT tagger

Training corpus \	WSJ	GENIA
WSJ	97.0	84.3
GENIA	75.2	98.1
<b>WSJ+GENIA</b>	<b>96.9</b>	<b>98.1</b>

Training corpus \	WSJ	GENIA
WSJ	96.7	84.3
GENIA	80.1	97.9
WSJ+GENIA	96.5	97.5

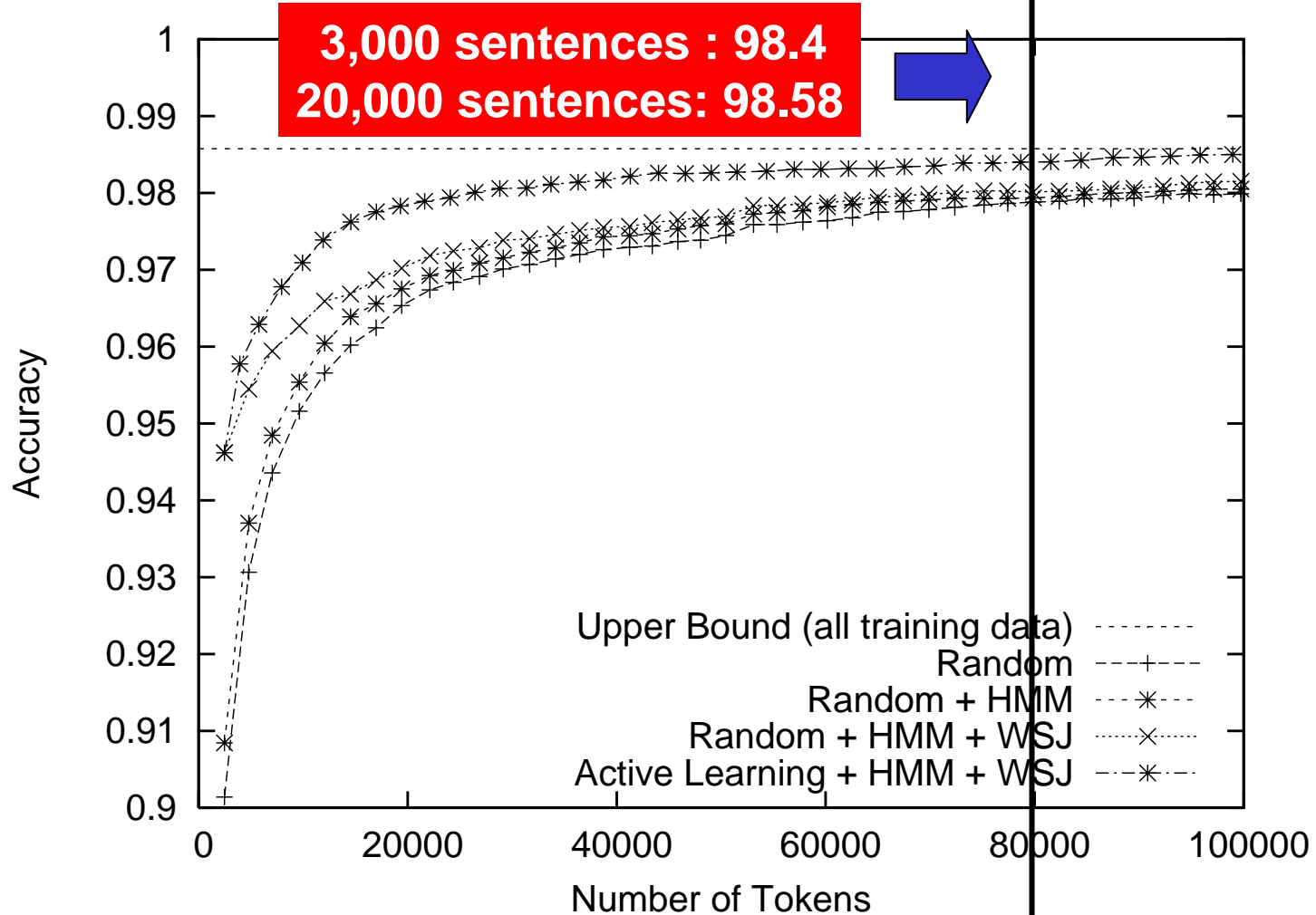
No degradation of the tagger  
trained by the mixed corpus

Some degradations (0.2 ~ 0.4)  
were observed, compared with  
the taggers trained by “pure”  
corpora



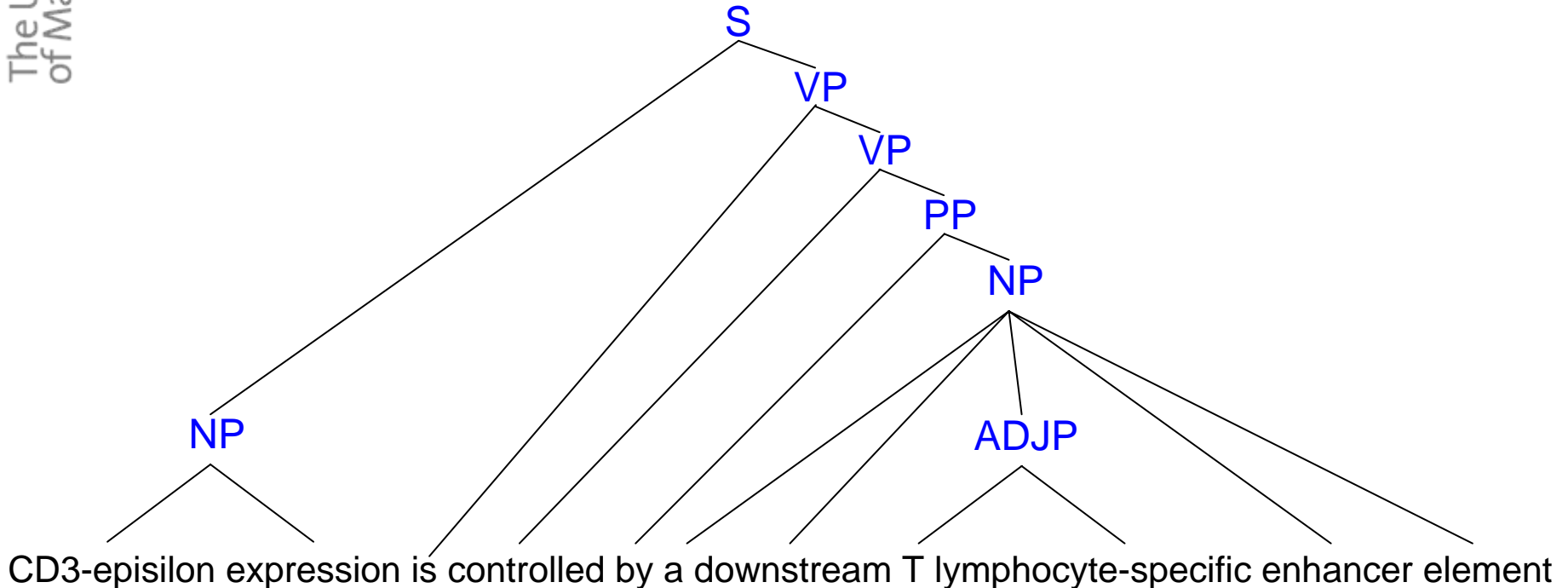
# CRF-based POS + Active Learning

## GENIA



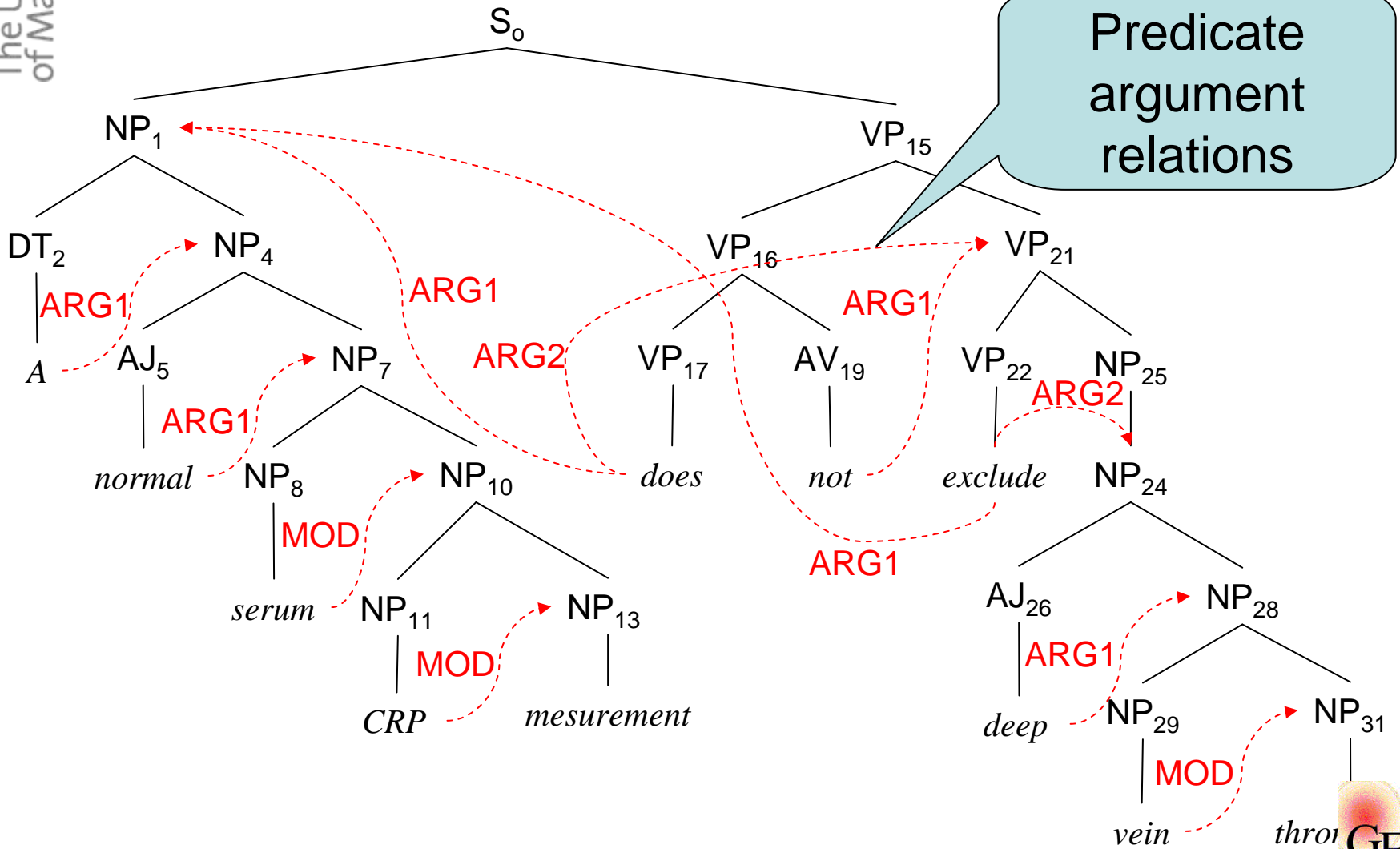
# Tool2 Full & Deep Parser Enju

# GENIA treebank



- Based on the standard of the Penn Tree Bank

# Semantic structure





# Performance of Semantic Parser

	Penn Treebank	GENIA
Coverage	99.7%	99.2%
F-Value (PA relations)	87.4%	86.4%
Sentence Precision	39.2%	31.8%
Processing Time	0.68sec	1.00sec

# Adaptation with Reference Distribution

**Lexical Assignment**

**Syntactic Preference**

$$p_E(t | \mathbf{w}) = \frac{1}{Z_{\mathbf{w}}} \prod_{w_i \in \mathbf{w}} p_{lex}(l_i | w_i) \cdot q_{syn}(t | \mathbf{I}),$$

$$Z_{\mathbf{w}} = \sum_{t \in T(\mathbf{w})} \prod_{w_i \in \mathbf{w}} p_{lex}(l_i | w_i) \cdot q_{syn}(t | \mathbf{I})$$

$$p_M(t | s) = \frac{1}{Z'_s} p_0(t | s) \exp \left( \sum_j \rho_j g_j(t | s) \right)$$

Feature function

Feature weight

Original model

	F-score	Training Time ( Sec )
Baseline ( PTB-trained, PTB-applied)	89.81	0
Baseline (PTB-trained, GENIA-applied)	86.39	0
Retraining ( GENIA )	88.45	14,695
Retraining ( PTB+GENIA )	89.94	238,576
Structure with RefDist	88.18	21,833
Lexical with RefDist	89.04	12,957
Lex/Structure with RefDist	90.15	31,637

	P	R	F	$\sigma_F$	AUC	$\sigma_{AUC}$
without evaluating self	60.4	69.3	64.3	4.3	0.879	0.026
without predicting self	60.4	65.6	62.7	3.5	0.834	0.032
with the prediction of self	57.8	66.1	61.4	3.9	0.914	0.020
(Miyao et al., 2008)	54.9	65.5	59.5			
(Airola et al., 2008)	52.9	61.8	56.4	5.0	0.848	0.023
(Sætre et al., 2007)	64.3	44.1	52.0			
(Mitsumori et al., 2006)	54.2	42.6	47.7			
(Yakushiji et al., 2005)	33.7	33.1	33.4			

Table 5: Comparison with previous results of the PPI extraction methods with the abstract-wise 10-fold cross validation on the AImed corpus

# Event Annotation

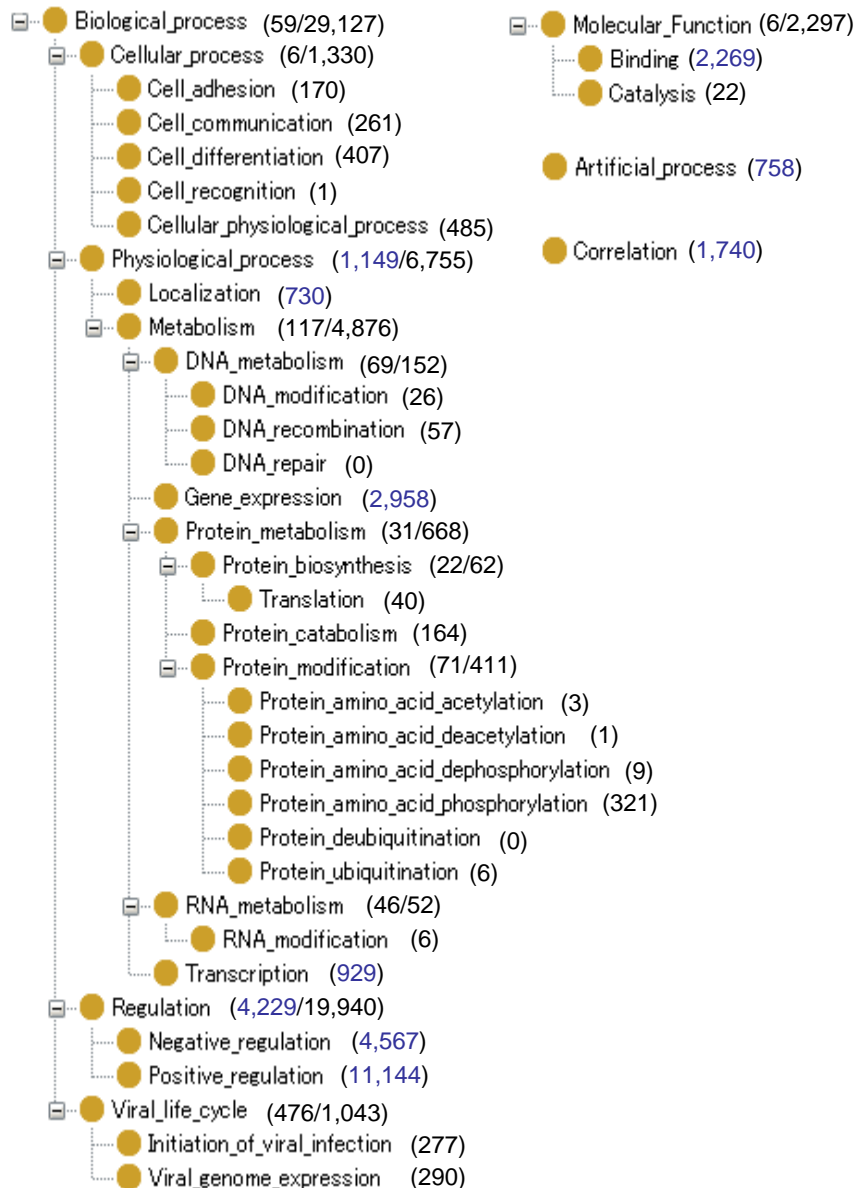
# GENIA event annotation

- Target of GENIA event annotation
  - Corpus
    - Part of GENIA corpus which is taken from PubMed using the MeSH terms, Human, Blood Cells and Transcription Factors.
  - Ontology
    - From the Gene ontology, concepts required for describing NF $\kappa$ B pathway have been selected (34 terms).
    - 3 additional concepts have been defined
      - Gene\_expression
      - Artificial\_process
      - Correlation

# GENIA event annotation – Stat (1/2)

- Annotation
  - 5 annotators + 1 manager with biology background.
  - using XConc Annotation tool
- 1,000 abstracts have been annotated
  - # of sentences: 8,981
  - # of sentences with events: 8,265
    - 92.0%
  - # of events: 34,065
    - Avg. 4.15 events/sentence

# GENIA event annotation – Stat (2/2)



- Correlation
  - meaning ‘some’ relation between events.
- Artificial\_process
  - Artificially performed processes.
  - Transfection, treatment, ...
- Gene\_expression
  - Transcription + Translation



# GENIA Event Annotation - example

Secretion of TNF<T32, the product of another NF-kappa B<T34-dependent gene<T33, was abolished by BHA<T35 in PMA-stimulated U937 cells<T37<T36.

## EVENT E23

TYPE : Localization

THEME : T32

CLUE : Secretion of TNF, the product of another NF-kappa B-dependent gene, was abolished by BHA in PMA-stimulated U937 cells.

ClueType

LinkThem

## EVENT E24

TYPE : Negative\_regulation

THEME : E23

CAUSE : T35

CLUE : Secretion of TNF, the product of another NF-kappa B-dependent gene, was abolished by BHA in PMA-stimulated U937 cells.

ClueLoc

LinkCause

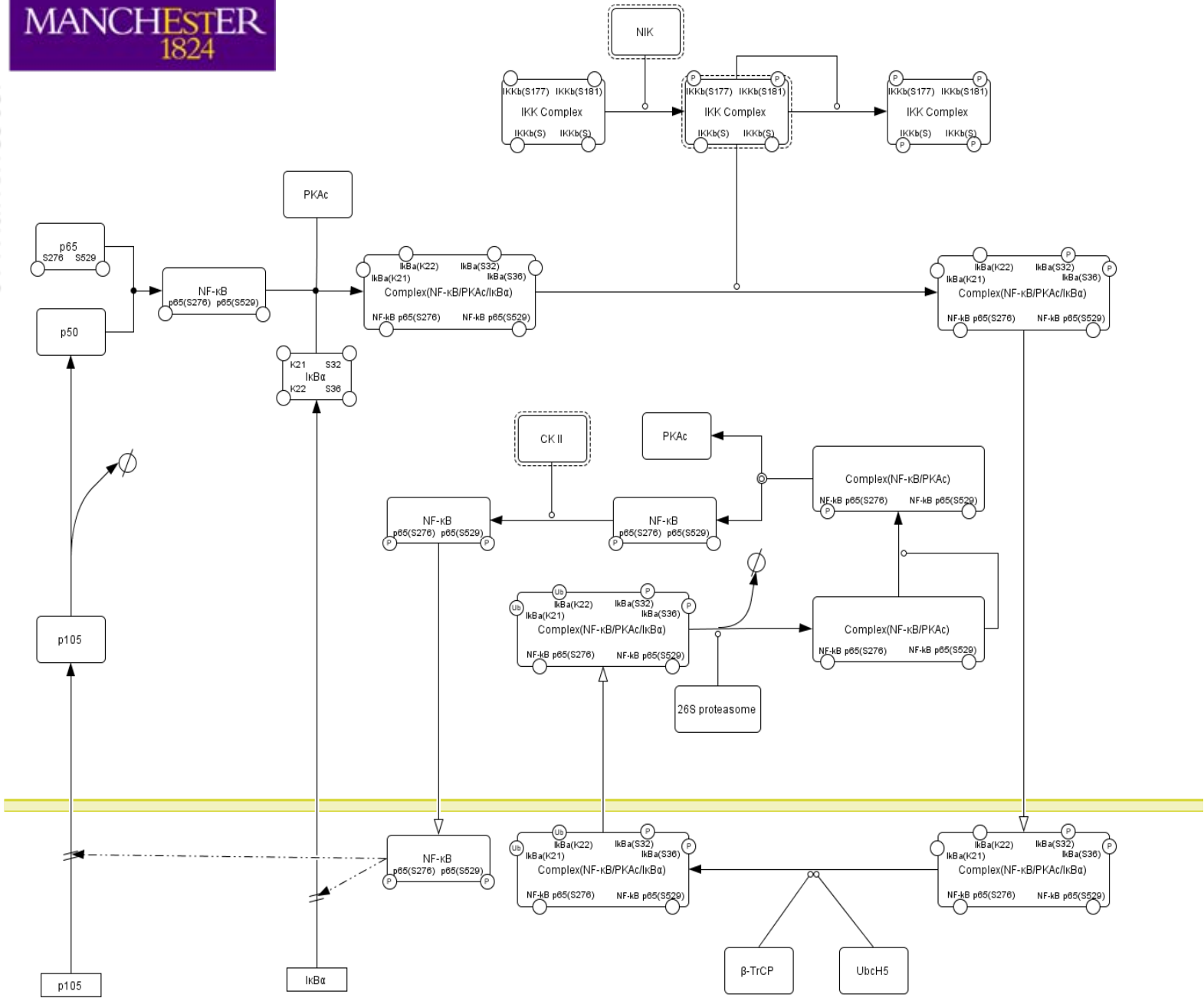
ClueType

- For an identified event in the given sentence,
  - classify the **type** of events and record the text span giving the clue of it (ClueType).
  - identify the **theme** of the events and record the text span linking the theme to the event (LinkTheme).
  - identify the **cause** of the events and record the text span linking the cause to the event (LinkCause).
  - record the environment (location, time) of the events (ClueLoc, ClueTime).

# Gene\_expression

- Theme patterns observed (2,958)
  - Protein 2,308
  - DNA 591
  - RNA 25
  - Peptide 4
  - Protein Protein 2
  - Erroneous 27
- Keywords
  - coexpress, nonexpress, overexpress, express, biosynthesis, product, synthesize, constitute, ...

# **Metabolic Pathways and Text Annotation**



# Discrepancy between event annotation and pathway annotation

<b>React.</b>	<b>Type</b>	<b>expected annotation class</b>	<b>○</b>	<b>×</b>
<b>R1</b>	<b>Binding</b>	<b>Binding</b>	<b>11</b>	<b>56</b>
<b>R2</b>	<b>Binding</b>	<b>Binding</b>	<b>16</b>	<b>20</b>
<b>R3</b>	<b>Phosphorylation</b>	<b>Protein_amino_acid_phosphorylation</b>	<b>19</b>	<b>1</b>
	<b>+ site information</b>	<b>Protein_amino_acid_phosphorylation</b>	<b>2</b>	<b>5</b>
	<b>+ modifier information</b>	<b>+ (any) Regulation</b>	<b>4</b>	<b>1</b>
<b>R4</b>	<b>ubiquitination</b>	<b>Protein_ubiquitination</b>	<b>5</b>	<b>0</b>
<b>R5</b>	<b>degradation</b>	<b>Protein_catabolism</b>	<b>27</b>	<b>1</b>
	<b>+ modifier information</b>	<b>+ (any) Regulation</b>	<b>6</b>	<b>3</b>
<b>R6</b>	<b>translocation</b>	<b>Localization</b>	<b>26</b>	<b>0</b>
<b>R7</b>	<b>Binding</b>	<b>Binding</b>	<b>23</b>	<b>12</b>
<b>R8</b>	<b>gene expression</b>	<b>Gene_expression (or Transcription + Translation)</b>	<b>5</b>	<b>9</b>
<b>R9</b>	<b>gene expression</b>	<b>Gene_expression (or Transcription + Translation)</b>	<b>1</b>	<b>10</b>
<b>R10</b>	<b>processing</b>	<b>Protein_processing</b>	<b>5</b>	<b>12</b>
		<b>total</b>	<b>150</b>	<b>130</b>

# Recommendations

- Think globally and act locally:
  - Semantics/Pragmatics: Domain and Task Specific
  - Processing/Application Guided Approach
- Tools Development
  - Tools based on standard
  - Adaptation methods for different domains
  - Annotation Tools for active learning
- Theory guided standardization
  - Semantics/Pragmatics: Theories on un-observables
  - Non-trivial Ontologies

# Tools and Corpora Available at

- National Centre for Text Mining
  - <http://www.nactem.ac.uk/>
- Tsujii Lab, Univ. of Tokyo
  - <http://www-tsujii.is.s.u-tokyo.ac.jp/>
- TerMine, AcroMine, AcroDisamb, GENIA Tagger, Shallow Parser, Deep Parser, Named Entity Recognizers, Event Recognizer (to be released),
- MEDIE, Info-Pubmed
- Interoperable Software Platforms (UIMA)
  - Contact: [kano@is.s.u-tokyo.ac.jp](mailto:kano@is.s.u-tokyo.ac.jp)