# Module 2: Introduction to IE and ANNIE

# About this tutorial

This tutorial comprises the following topics:

- Introduction to IE

- ANNIE

- Multilingual tools in GATE

- Evaluation and Corpus Quality Assurance

In Module 3, you'll learn how to use JAPE, the pattern matching language that many PRs use
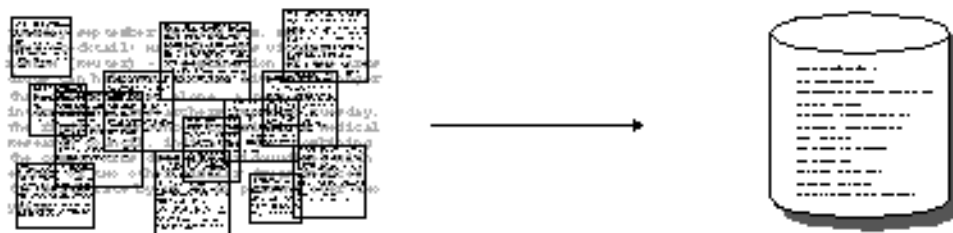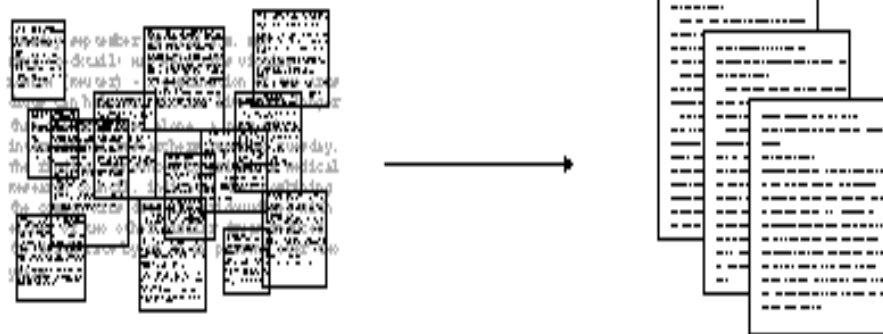
# What is information extraction?

# IE is not IR

- IR pulls **documents** from large text collections (usually the Web) in response to specific keywords or queries. You analyse the **documents**.

- IE pulls **facts** and **structured information** from the content of large text collections. You analyse the **facts**.

# IE for Document Access

- With traditional query engines, getting the facts can be hard and slow

    - Where has the Queen visited in the last year?

    - Which airports are currently closed due to the volcanic ash?

- Which search terms would you use to get thess?

- How can you specify you want to see someone's home page?

- IE returns information in a structured way

- IR returns documents containing the relevant information somewhere

# IE as an alternative to IR

- IE returns knowledge at a much deeper level than traditional IR

- It allows you to specify your query in a more structured way

- Constructing a database through IE and linking it back to the documents can provide a valuable alternative search tool

- Even if results are not always accurate, they can be valuable if linked back to the original text

# What is IE used for?

- IE is an enabling technology for many other applications:
  - Text Mining
  - Semantic Annotation
  - Question Answering
  - Opinion Mining
  - Decision Support
  - Rich information retrieval and exploration
  - and so on..

# Two main types of IE systems

## Knowledge Engineering

- rule based
- developed by experienced language engineers
- make use of human intuition
- require only small amount of training data
- development can be very time consuming
- some changes may be hard to accommodate

## Learning Systems

- use statistics or other machine learning
- developers do not need as much LE expertise
- require large amounts of annotated training data
- some changes may require re-annotation of the entire training corpus

# Named Entity Recognition: the cornerstone of IE

Traditionally, NER is the identification of proper names in texts, and their classification into a set of predefined categories of interest

- Person

- Organisation (companies, government organisations, committees, etc)

- Location (cities, countries, rivers, etc)

- Date and time expressions

Various other types are frequently added, as appropriate to the application, e.g. newspapers, ships, monetary amounts, percentages etc.

# Why is NER important?

- NER provides a foundation for building more complex IE systems

- Relations between NEs can provide tracking, ontological information and scenario building

- Tracking (co-reference): "Dr Smith", "John Smith", "John", "he"

- Ontologies: "Athens, Georgia" vs "Athens, Greece"

# Typical NE pipeline

- Pre-processing (tokenisation, sentence splitting, morphological analysis, POS tagging)

- Entity finding (gazetteer lookup, NE grammars)

- Coreference (alias finding, orthographic coreference etc.)

- Export to database / XML / ontology

# Example of IE

John lives in London. He works there for Polar Bear Design.

# Basic NE Recognition

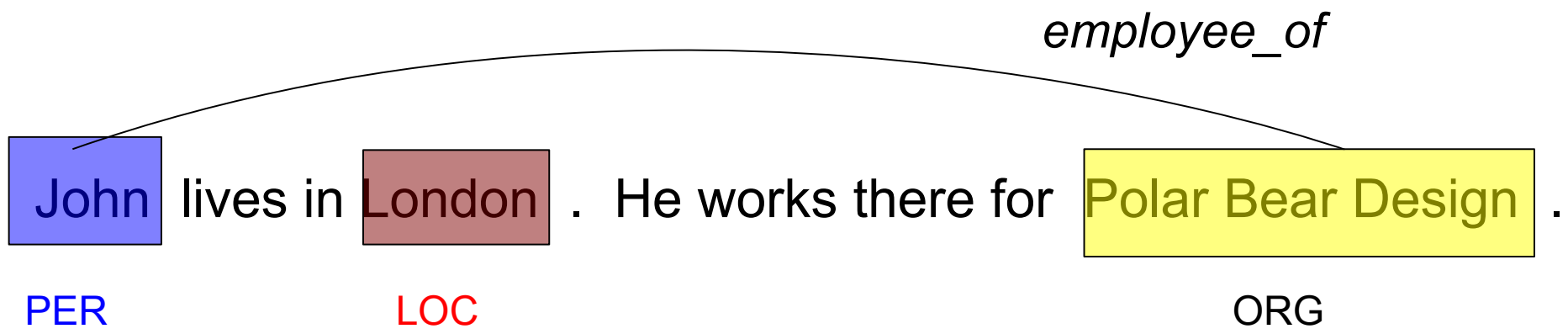John lives in London . He works there for Polar Bear Design .

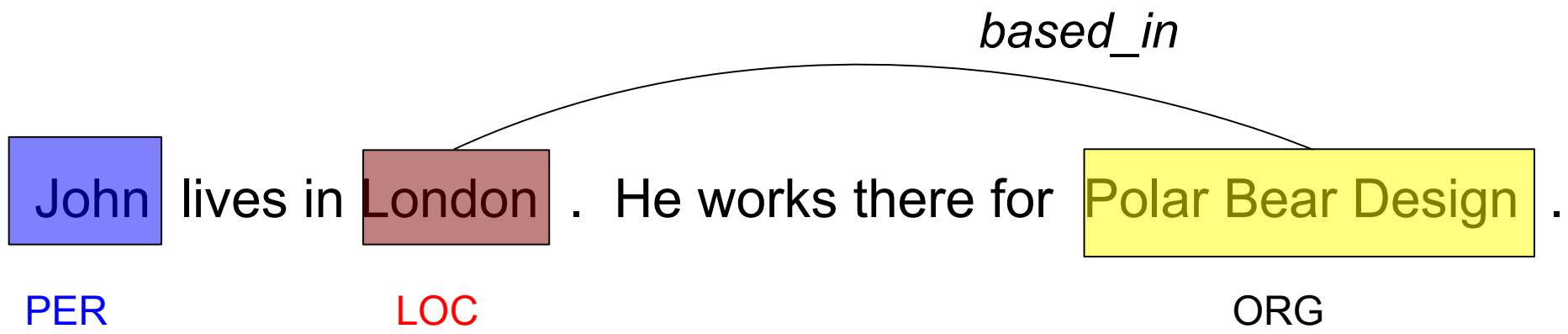PERSON          LOCATION                          ORGANISATION

# Co-reference

*same_as*

John lives in London . He works there for Polar Bear Design .

PER        LOC                              ORG

# Relations

*live_in*

John lives in London . He works there for Polar Bear Design .

PER          LOC                                        ORG

# Relations (2)

*employee_of*

John lives in London . He works there for Polar Bear Design .

PER       LOC       ORG

# Relations (3)

*based_in*

John lives in London .  He works there for Polar Bear Design .

PER           LOC                      ORG

# **Examples of IE systems**

# HaSIE

- Health and Safety Information Extraction

- Application developed with GATE, which aims to find out how companies report about health and safety information

- Answers questions such as:

  - "How many members of staff died or had accidents in the last year?"

  - "Is there anyone responsible for health and safety?"

- IR returns whole documents

## Hse

**CompanyName**

BAA

**HSEParagraphs**

sustainability management system. ... BAA has received a RoSPA gold award for occupational safety for the fourth year running. The award is given only if a consistently good or continuously improving performance can be demonstrated over a four-year period. The accident frequency ratio for construction projects was 0.4 (0.49) per 100,000 hours worked, less than one third of the national accident frequency rate in the construction sector. The company is running a ?One in a Million? campaign to raise safety consciousness and standards in construction and reduce the accident frequency rate still further to one for every million man hours worked. ... We have no higher priority than the safety and security of the passengers, staff and organisations that use our airports. In order to ensure that our systems and practices are continually assessed and upgraded, we work

**Awards**

BAA has received a RoSPA gold award

**Accidents**

The accident frequency ratio for construction projects was 0.4 (0.49) per 100,000 hours worked, less than one third of the national accident frequency rate in the construction sector.

# Obstetrics records

- Streamed entity recognition during note taking

  - Interventions, investigations, etc.

- Based entirely on gazetteers and JAPE

- Has to cope with terse, ambiguous text and distinguish past events from present

- Used upstream for decision support and warnings

**File  Options  Tools  Help**

GATE

- Applications
  - pipeline
- Language Resources
- Processing Resources
  - Cleanup
  - Annotation Set Tra...
  - IE Transducer
  - Flexible Gazetteer
  - Roots gazetteer

| MimeType | ▼ | text/l... |
| currentGravidity | ▼ | 3 |
| day | ▼ | 20 |
| gate.SourceURL | ▼ | file:/l... |
| month | ▼ | 8 |
| shift | ▼ | 12 |

Rename this resource

Messages | pipeline | Case_006.htm_00...

**Annotation Sets  Annotations List  Co-reference Editor  Text**

1:30pm

Cx: 3cm. contractions q2–3min.   FHR: reassuring.   reactive.

4:00pm

BP: 140/90.

PV: 6cm; 60%; −1; soft consistency, anterior position; cephalic; Intact membranes; no vaginal bleeding.

Contractions: 3/10min; regular; moderate

On urinalysis: Protein > 300mg

BP before 20 weeks gestation: 120/80

Plan: monitor Vital Signs by protocol for elevated BP

5:15pm

| Type | Set | St |

18 Annotations (0 selected)  Select:

**Document Editor  Initialisation Parameters**

- ☐ CesareanSectionInPriorDelivery
- ☑ DiastolicBloodPressure
- ☑ DiastolicBloodPressureBefore20W
- ☑ Dinoprostone
- ☐ EstimatedFetalWeight
- ☐ FHREvaluation
- ☐ GBSNeonatalSepsisAfterAPrevious
- ☐ Gravidity
- ☐ HighRiskForAnaphylaxis
- ☐ MagnesiumSulfate
- ☑ MembranesStatus
- ☐ MyastheniaGravis
- ☐ PatientAge
- ☐ PelvicAdequacy
- ☐ PenicillinAllergy
- ☐ PreviousCesareanSectionType
- ☑ SystolicBloodPressure
- ☑ SystolicBloodPressureBefore20We
- ☐ TimeStamp
- ☑ UrineProtein

New

# Multiflora

- IE system in the botanical domain

- Finds information about different plants: size, leaf span, colour etc

- Collates information from different sources: these often refer to plant features in slightly different ways

- Uses shallow linguistic analysis: POS tags and noun and verb phrase chunking

- Important to relate features to the right part of the plant: leaf size rather than plant size, colour of flowers vs colour of leaves etc.

Gate 2.1-beta1 build 1061

File  Options  Tools  Help

**Messages** | 📄 **R_a_FNA.txt_00743**

**Text** | **Annotations** | **Annotation Sets** | **Print** | 🖨

7. Ranunculus acris Linnaeus, Sp. Pl. 1: 554. 1753
☐ Renoncule âcre, bouton d'or
Ranunculus acris var. latisectus Beck
Stems erect from short caudex or rhizome, never rooting nodally, hispid,
strigose, or glabrous, base not bulbous. Roots never tuberous. Basal leaf
blades pentagonal in outline, deeply 3–5–parted, 1.8–5.2 X 2.7–9.8 cm,
segments 1–2 X –lobed or –parted, ultimate segments narrowly elliptic or
oblong to lanceolate, margins toothed or lobulate, apex acute to rounded.
Flowers: receptacle glabrous; sepals spreading, 4–6(–9) x 2–5 mm, hispid;
petals 5, yellow, 8–11(–17) X 7–13 mm. Heads of achenes globose,
5–7(–10) mm wide; achenes 2–3 X 1.8–2.4 mm, glabrous, margin forming
narrow rib 0.1–0.2 mm wide; beak persistent, deltate, usually with tip short
or long, straight or curved, subulate, 0.2–1 mm. 2n = 14.

| Type | Set | Start ▲ | End | Features |
|------|-----|---------|-----|----------|
| PlantFeatures | Default | 0 | 1 | {type=number} |
| Header | Default | 0 | 44 | {} |
| PlantFeatures | Default | 38 | 39 | {type=number} |
| PlantFeatures | Default | 103 | 113 | {rule=HeadAdj} |
| Head | Default | 119 | 124 | {} |
| PlantFeatures | Default | 125 | 130 | {rule=HeadAdj} |
| PlantFeatures | Default | 136 | 141 | {rule=AdjHead} |
| Head | Default | 142 | 148 | {} |
| Head | Default | 152 | 159 | {} |

**Annotations Editor** | **Features Editor**

Default annotations
- ☑ Head
- ☑ Header
- ☐ Lookup
- ☑ PlantFeature
- ☐ Segment
- ☐ SegmentSplit
- ☐ Sentence
- ☐ SpaceToken
- ☐ Split
- ☐ Token

Original markups a
- ☐ paragraph

Hides this view

# Old Bailey IE

- The Old Bailey Proceedings Online makes available a fully searchable, digitised collection of all surviving editions of the Old Bailey Proceedings from 1674 to 1913

- GATE was used to perform IE on the court reports, identifying names of people, places, dates etc.

- ANNIE was customised to only extract full Person names and to take account of old English language used

- More info at http://www.oldbaileyonline.org/static/Project.jsp

# Old Bailey IE

# IE in other languages

- ANNIE has been adapted to various other languages: some as test cases, some as real IE systems

- More details about this in Track 3 (Advanced IE module)

- Brief introduction to multilingual PRs in GATE later in this tutorial

Gate 2.1_02-beta build 1299

File   Options   Tools   Help

te
**Applications**
  arabic not trained
  **Language Resources**
  GATE document_00095
  **Processing Resources**
  orthomatcher
  arabic not trained grammar
  arabic gaz
  arabic tokeniser
  reset
  **Data stores**
  file:/share/nlp.18/diana/ga

loaded in 2.677 seconds

GATE document_00095
file:/share/nlp.18/diana/gatecorpora/arabic/treebank/bbnfiles/test/processed/
Messages

Text   Annotations   Annotation Sets   Print

Default annotations
Key annotations
  ☑ Cardinal
  ☑ Date
  ☑ Event
  ☑ Gpe
  ☑ Gpe_desc
  ☑ Money
  ☑ Nationality
  ☑ Ordinal
  ☑ Org_desc
  ☑ Organization
  ☑ Per_desc
  ☑ Person
Original markups annota

Annotations Editor   Features Editor   Initialisation Parameters

File   Options   Tools   Help

Messages   📄 BengaliSampleText.utf8.txt

আমার নাম অনিল রায়। আমা ল্যন্কাস্টরে থাকা। আমার বাবা লিওয়ারপুলে থাকে।

আমার বাবার নাম হচ্ছে রাজেশে রায়। ল্যন্কাস্টর ইউনিওয়্যরসিটি আমার পদার যায়গা। আমার বাবা কোকা কোলা কম্পনিতে কাজ করে।

My name is Anil Roy. I live in Lancaster. My father lives in Liverpool. My

father's name is Rajesh Roy. Lancaster University is my place of

**Default annotations**

- [ ] DEFAULT_TOKEN
- [x] Location
- [ ] Lookup
- [x] Organisation
- [x] Person
- [ ] SpaceToken
- [ ] Token

| Type | Set | Start ▲ | End | Features |
|------|-----|---------|-----|----------|
| Person | Default | 10 | 18 | {kind=fullname} |
| Location | Default | 27 | 38 | {kind=city, rule=City} |
| Location | Default | 59 | 67 | {kind=city, rule=City} |
| Person | Default | 101 | 112 | {kind=fullname} |
| Organisation | Default | 115 | 141 | {} |
| Organisation | Default | 173 | 182 | {} |

Annotations   Features

Gate

- Applications
  - Bengali NE
- Language Resources
  - BengaliSampleText.utf8.t
- Processing Resources
  - BengaliNE
  - BengaliTokeniser
  - bengali_gazetteer
- Data stores

Bengali NE run in 0.591 seconds

# ANNIE: A Nearly New Information Extraction system

# Nearly New Information Extraction

- ANNIE is a readymade collection of PRs that performs IE on unstructured text.

- For those who grew up in the UK, you can think of it as a Blue Peter-style "here's one we made earlier".

- ANNIE is "nearly new" because

  - It was based on an existing IE system, LaSIE

  - We rebuilt LaSIE because we decided that people are better than dogs at IE

  - Being 10 years old, it's not really new any more

# What's in ANNIE?

- The ANNIE application contains a set of core PRs:

  - Tokeniser

  - Sentence Splitter

  - POS tagger

  - Gazetteers

  - Named entity tagger (JAPE transducer)

  - Orthomatcher (orthographic coreference)

- There are also other PRs  available in the ANNIE plugin, which are not used in the default application, but can be added if necessary

  - NP and VP chunker

# Core ANNIE components

# Loading and running ANNIE

- Because ANNIE is a ready-made application, we can just load it directly from the menu

- Click the      icon from the top GATE menu OR Select File → Load ANNIE system

- Select "with defaults"

- Load any document from <u>module-2-hands-on/news-texts</u> and add it to a corpus

- Run ANNIE and inspect the annotations

- You should see a mixture of Named Entity annotations (Person, Location etc) and some other linguistic annotations (Token, Sentence etc)

# Let's look at the PRs

- Each PR in the ANNIE pipeline creates some new annotations, or modifies existing ones

- Document Reset → removes annotations

- Tokeniser → Token annotations

- Sentence Splitter → Sentence, Split annotations

- Gazetteer → Lookup annotations

- POS tagger → adds category features to Token annotations

- JAPE transducer → Date, Person, Location, Organisation, Money, Percent annotations

- Orthomatcher → adds match features to NE annotations

# Document Reset

- This PR should go at the beginning of (almost) every application you create

- It removes annotations created previously, to prevent duplication if you run an application more than once

- By default it does not remove the "Original markups" AS

- You can configure it to keep any other annotation sets you want, or to remove particular annotation types only

# Document Reset Parameters



Specify any specific annotations to remove. By default, remove all.

Keep Original Markups set

Keep Key set

# Document Reset Parameters



Specify any specific annotations to remove. By default, remove all.

Remove the Temp set only

# Tokenisation and sentence splitting

# **Tokeniser**

- Tokenisation based on Unicode  classes

- Declarative token specification language

- Produces Token and SpaceToken annotations with features orthography and kind

- Length and string features are also produced

- Rule for a lowercase word with initial uppercase letter

```
"UPPERCASE_LETTER" LOWERCASE_LETTER"* >
  Token; orthography=upperInitial; kind=word
```

# Document with Tokens

# ANNIE English Tokeniser

- The English Tokeniser is a slightly enhanced version of the Unicode tokeniser

- It comprises an additional JAPE transducer which adapts the generic tokeniser output for the POS tagger requirements

- It converts constructs involving apostrophes into more sensible combinations (like the Penn Treebank)

  - don't  →  do + n't

  - you've → you + 've

# **Looking at Tokens**

- Tidy up GATE by removing all resources and applications (or just restart GATE)

- Create a new corpus and populate it with module

- Create a new application (corpus pipeline)

- Load a Document Reset and an ANNIE English Tokeniser

- Add them (in that order) to the application

- Set the Document Reset to keep the Key AS

- Run the application on the corpus

- View the Token and SpaceToken annotations

- What different values of the "kind" feature do you see?

43

# Looking at Tokens



House prices in England and Wales were 10.8 per cent higher in the second quarter than in the same period last year, according to figures released by the Land Registry.

The average house cost almost £7,000 ($10,000) more than it did in the same period of 2000.

| Type | Set | Start | End | Id | Features |
|------|-----|-------|-----|-----|----------|
| Token | | 0 | 5 | 97 | {kind=word, length=5, orth=upperInitial, string=House} |
| Token | | 6 | 12 | 99 | {kind=word, length=6, orth=lowercase, string=prices} |
| Token | | 13 | 15 | 101 | {kind=word, length=2, orth=lowercase, string=in} |
| Token | | 16 | 23 | 103 | {kind=word, length=7, orth=upperInitial, string=England} |
| Token | | 24 | 27 | 105 | {kind=word, length=3, orth=lowercase, string=and} |
| Token | | 28 | 33 | 107 | {kind=word, length=5, orth=upperInitial, string=Wales} |
| Token | | 34 | 38 | 109 | {kind=word, length=4, orth=lowercase, string=were} |
| Token | | 39 | 41 | 111 | {kind=number, length=2, string=10} |
| Token | | 41 | 42 | 112 | {kind=punctuation, length=1, string=.} |
| Token | | 42 | 43 | 113 | {kind=number, length=1, string=8} |
| Token | | 44 | 47 | 115 | {kind=word, length=3, orth=lowercase, string=per} |
| Token | | 48 | 52 | 117 | {kind=word, length=4, orth=lowercase, string=cent} |
| Token | | 53 | 59 | 119 | {kind=word, length=6, orth=lowercase, string=higher} |
| Token | | 60 | 62 | 121 | {kind=word, length=2, orth=lowercase, string=in} |

# Sentence Splitter

- The default splitter finds sentences based on Tokens

- Creates Sentence annotations and Split annotations on the sentence delimiters

- Uses a gazetteer of abbreviations etc. and a set of JAPE grammars which find sentence delimiters and then annotate sentences and splits

- Load an ANNIE Sentence Splitter and add it to your application (after the tokeniser)

- Run the application and view the results

# Document with Sentences

# Sentence splitter variants

- An alternate set of rules can be loaded with the regular sentence splitter

- To do this, load "main-single-nl.jape" instead of "main.jape" as the value of the grammar parameter

- The main difference is the way it handles new lines

- In some cases, you might want a new line to signal a new sentence, e.g. addresses

- In other cases, you might not, e.g. in emails that have been split by the email program

- A regular expression Java-based splitter is also available, called RegEx Sentence Splitter, which is sometimes faster

- This handles new lines in the same way as the default sentence splitter

- See "Further Exercises" to experiment with splitter variants

University of Sheffield NLP

# **Shallow lexico-syntactic features**

# POS tagger

- ANNIE POS tagger is a Java implementation of Brill's transformation based tagger

- Previously known as **Hepple Tagger** (you may find references to this and to **heptag**)

- Trained on the Wall Street Journal corpus (news articles), uses Penn Treebank tagset

- Default ruleset and lexicon can be modified manually (with a little deciphering)

- Adds a <u>category</u> feature to Token annotations

- Requires Tokeniser and Sentence Splitter to be run first

# Morphological analyser

- This is in the Tools plugin (not an integral part of ANNIE)

- Flex based rules: can be modified by the user (instructions in the User Guide)

- Generates "root" feature on Token annotations

- The tokeniser must be run before the morpher

- The POS tagger must run before the morpher if the *considerPOSTag* parameter is set to true (and the POS tagger requires a sentence splitter first)

# Shallow lexico-syntactic features

- Add an ANNIE POS Tagger to your app

- Add a GATE Morphological Analyser after the POS Tagger

- If this PR is not available, load the Tools plugin first

- Examine the features of the Token annotations

  - New features of category and root have been added

# Shallow lexico-syntactic features

House prices in England and Wales were 10.8 per cent higher in the second quarter than in the same period last year, according to figures released by the Land Registry.

The average house cost almost £7,000 ($10,000) more than it did in the same period of 2000.

Prices rose in all regions, with the East Midlands up the most at 15.6 per cent, while house prices in Wales experienced the smallest increase at 6.5 per cent.

| Type | Set | Start | End | Id | Features |
|------|-----|-------|-----|-----|----------|
| Token | | 0 | 5 | 609 | {category=NN, kind=word, length=5, orth=upperInitial, root=house, string=H |
| Token | | 6 | 12 | 611 | {affix=s, category=NNS, kind=word, length=6, orth=lowercase, root=price, s |
| Token | | 13 | 15 | 613 | {category=IN, kind=word, length=2, orth=lowercase, root=in, string=in} |
| Token | | 16 | 23 | 615 | {category=NNP, kind=word, length=7, orth=upperInitial, root=england, strin |
| Token | | 24 | 27 | 617 | {category=CC, kind=word, length=3, orth=lowercase, root=and, string=and} |
| Token | | 28 | 33 | 619 | {category=NNP, kind=word, length=5, orth=upperInitial, root=wales, string= |
| Token | | 34 | 38 | 621 | {affix=ed, category=VBD, kind=word, length=4, orth=lowercase, root=be, st |
| Token | | 39 | 41 | 623 | {category=CD, kind=number, length=2, root=10, string=10} |
| Token | | 41 | 42 | 624 | {category=., kind=punctuation, length=1, root=., string=.} |
| Token | | 42 | 43 | 625 | {category=CD, kind=number, length=1, root=8, string=8} |
| Token | | 44 | 47 | 627 | {category=IN, kind=word, length=3, orth=lowercase, root=per, string=per} |
| Token | | 48 | 52 | 629 | {category=NN, kind=word, length=4, orth=lowercase, root=cent, string=cen |
| Token | | 53 | 59 | 631 | {category=JJR, kind=word, length=6, orth=lowercase, root=higher, string=hi |
| Token | | 60 | 62 | 633 | {category=IN, kind=word, length=2, orth=lowercase, root=in, string=in} |

# Gazetteers

# Gazetteers

- Gazetteers are plain text files containing lists of names (e.g rivers, cities, people, …)

- The lists are compiled into Finite State Machines

- Each gazetteer has an index file listing all the lists, plus features of each list (majorType, minorType and language)

- Lists can be modified either internally using the Gazetteer Editor, or externally in your favourite editor (note that the new Gazetteer editor replaces the old GAZE editor you may have seen previously)

- Gazetteers generate Lookup annotations with relevant features corresponding to the list matched

- Lookup annotations are used primarily by the NE transducer

- Various different kinds of gazetteer are available: first we'll look at the default ANNIE gazetteer

# Gazetteer editor



definition file entries

entries for selected list

# ANNIE gazetteer

- Create an ANNIE Gazetteer PR, but change the listsURL to gazetteer/lists.def in the hands-on materials for this module. Double-click on it to open.

- Select "Gazetteer Editor" from the bottom tab

- In the left hand pane (linear definition) you see the index file containing all the lists

- In the right hand pane you see the contents of the list selected in the left hand pane

- Each entry can be edited by clicking in the box and typing

- New lists and entries can be added by typing in the "New list" or "New entry" box respectively

# Modifying the definition file

add a new list

edit an existing list name by typing here

edit the major and minor Types by typing here

delete a list by right clicking on an entry and selecting Delete

| airport.lst ▼ | New List | | |
|---|---|---|---|
| **List name** | **Major** | **Minor** | **Language** |
| charities.lst | organization | | |
| city.lst | location | city | |
| city_cap.lst | location | city | |
| company.lst | organization | company | |
| company_cap.lst | organization | company | |
| country.lst | location | country | |
| country_abbrev.lst | location | country_abbrev | |
| country_adj.lst | country_adj | | |
| country_cap.lst | location | country | |
| currency_prefix.lst | currency_unit | pre_amount | |
| currency_unit.lst | currency_unit | post_amount | |
| date_key.lst | date_key | | |
| date_unit.lst | date_unit | | |
| day.lst | date | day | |
| day_cap.lst | date | day | |
| department.lst | organization | departmen | |

57

# Modifying a list

add a new entry
by typing here

edit an
existing entry
by typing here

Delete an entry by
right clicking and
selecting "Delete"

| | | |
|---|---|---|
| | New Entry | Add Cols |
| **Value** | | |
| Aaccra | | |
| Aalborg | | |
| Aarhus | | |
| Ababa | | |
| Abadan | | |
| Abakan | | |
| Aberdeen | | |
| Abha | | |
| Abi Dhabi | | |
| Abidjan | | |
| Abilene | | |
| Abu | | |
| Abu Dhabi | | |
| Abuja | | |
| Acapulco | | |
| Filter: | | 1993 entries |

58

# Editing gazetteer lists

- The ANNIE gazetteer has about 60,000 entries arranged in 80 lists

- Each list reflects a certain category, e.g. airports, cities, first names etc.

- List entries might be entities or parts of entities, or they may contain contextual information (e.g. job titles often indicate people)

- Click on any list to see the entries

- Note that some lists are not very complete!

- Try adding, deleting and editing existing lists, or the list definition file

- To save an edited gazetteer, right click on the gazetteer name in the tabs at the top or in the resources pane on the left, and select "Save and reinitialise". (You can select "Save as..." to save your changes in a different location.)

- Try adding a word from a document you have loaded (that is not currently recognised as a Lookup) into the gazetteer, re-run the gazetteer and check the results.

59

# Editing gazetteers outside GATE

- You can also edit both the definition file and the lists outside GATE, in your favourite text editor

- You need to reinitialise the gazetteer in GATE (in order to re-load the modified files) before running it again

- To reinitialise any PR, right click on its name in the Resources pane and select "Reinitialise"

# List attributes

- When something in the text matches a gazetteer entry, a Lookup annotation is created, with various features and values

- The ANNIE gazetteer has the following default feature types: majorType, minorType, language

- These features are used as a kind of classification of the lists: in the definition file features are separated by ":"

- For example, the "city" list has a majorType "location" and minorType "city", while the "country" list has "location" and "country" as its types

- Later, in the JAPE grammars, we can refer to all Lookups of type location, or we can be more specific and refer just to those of type "city" or type "country"

# NE transducers

# NE transducer

- Gazetteers can be used to find terms that suggest entities

- However, the entries can often be ambiguous

  - "May Jones" vs "May 2010" vs "May I be excused?"

  - "Mr Parkinson" vs "Parkinson's Disease"

  - "General Motors" vs. "General Smith"

- Handcrafted grammars are used to define patterns over the Lookups and other annotations

- These patterns can help disambiguate, and they can combine different annotations, e.g. Dates can be comprised of day + number + month

- Each NE transducer consists of one or more grammars written in the JAPE language, which Module 3 will cover in detail tomorrow

# ANNIE NE Transducer

- Create an ANNIE NE Transducer PR with the default parameters

- Add it to the end of the application

- Run the application

- Look at the annotations

- You should see some new annotations such as Person, Location, Date etc.

- These will have features showing more specific information (eg what kind of location it is) and the rules that were fired (for ease of debugging)

# Co-reference

# Using co-reference

- Different expressions may refer to the same entity

- Orthographic co-reference module (orthomatcher) matches proper names and their variants in a document

- [Mr Smith] and [John Smith] will be matched as the same person

- [International Business Machines Ltd.] will match [IBM]

# Orthomatcher PR

- Performs co-reference resolution based on orthographical information of entities

- Produces a list of annotation ids that form a co-reference chain

- List of such lists stored as a document feature named "MatchesAnnots"

- Improves results by assigning entity type to previously unclassified names, based on  relations with classified entities

- May not reclassify already classified entities

- Classification of unknown entities very useful for surnames which match a full name, or abbreviations,
  e.g. "Bonfield" <Unknown> will match "Sir Peter Bonfield" <Person>

- A pronominal co-reference PR is also available in the ANNIE plugin

# **Looking at co-reference**

- Add a new PR: ANNIE OrthoMatcher

- Add it to the end of the application

- Run the application

- In a document view, open the co-reference editor by clicking the button above the text

- All the documents in the corpus should have some co-reference, but some may have more than others

# Coreference editor

# Using the co-reference editor

- Select the annotation set you wish to view (Default)
- A list of all the co-reference chains that are based on annotations in the currently selected set is displayed
- Select an item in the list to highlight all the member annotations of that chain in the text (you can select more than one at once)
- Hovering over a highlighted annotation in the text enables you to Delete an item from the co-reference chain
- Try it!
- Deselect all items in this list, then select a type from the "Type" combo box and click "Show" to view all coreferences of a particular annotation type (note that some types may not have coreferences)

# Modifying ANNIE

# Modifying ANNIE

- Typically any new application you want to create will use some or all of the core components from ANNIE

- The tokeniser, sentence splitter and orthomatcher are basically language, domain and application-independent

- The POS tagger is language dependent but domain and application-independent

- You may also require additional PRs (either existing or new ones – e.g. morphological analyser

- The gazetteer lists and JAPE grammars may act as a starting point but will almost certainly need to be modified

# ANNIE without defaults

- This option loads all the ANNIE PRs, but enables you to change the location of any of them

- It's useful If you want to use ANNIE but you want to change some of the PRs slightly or replace them with your own modified versions

- Restart GATE or remove all PRs and applications, to tidy up a little

- Load ANNIE as before, but this time select "Without defaults"

- For each PR, select the default option, except for the gazetteer, where you should select your saved gazetteer index file (gazetteer/lists.def) from the hands-on materials

# Multilingual IE

# Language plugins

- Language plugins contain language-specific PRs, with varying degrees of sophistication and functions for:
  - Arabic
  - Cebuano
  - Chinese
  - Hindi
  - Romanian

- There are also various applications and PRs available for French, German and Italian

- These are not plugins, strictly speaking, as they do not provide new kinds of PRs

- Applications and individual PRs for these are found in gate/plugins directory: load them as any other PR

- More details of language plugins in user guide

# Building a language-specific application

- The following PRs are largely language-independent:
    - Unicode tokeniser
    - Sentence splitter
    - Gazetteer PR (but do localise the lists!)
    - Orthomatcher (depending on the nature of the language)
- Other PRs will need to be adapted (e.g. JAPE transducer) or replaced with a language-specific version (e.g. POS tagger)
- This topic is covered in more detail in Module 10 (Advanced IE, in Track 3)

# Useful Multilingual PRs

- Stemmer_Snowball plugin

  - Stemmer PRs with models for Danish, Dutch, English, Finnish, French, German, Italian, Norwegian, Portuguese, Russian, Spanish, Swedish (set the language init parameter with one of these in lower case)

  - Run a Tokeniser first (Unicode one is best)

- Tagger_Framework

  - GenericTagger PR which can run external tools such as the TreeTagger

# Displaying multilingual data

GATE uses standard (and imperfect) Java rendering engine for displaying text in multiple languages.

# Displaying multilingual data

All visualisation and editing tools use the same facilities

# Editing multilingual data

- Java currently provides no special support for text input, but the GATE Unicode Kit (GUK) allows the definition of additional Input Methods (currently 30 IMs for 17 languages)
- Pluggable in other applications (e.g. MPI's EUDICO)
- Can use virtual keyboard or standard layouts over QWERTY
- IMs defined in plain text files
- Use Options → Input methods to change the keyboard mapping
- **However**, if your OS provides a way to change keyboard mappings on the fly (as most do these days) and covers the language you want to use, you should use that function instead.

# Annotation and Evaluation

# Topics covered

- Defining annotation guidelines

- Recap on manual annotation using the GATE GUI

- Using the GATE evaluation tools

# Before you start annotating...

- You need to think about annotation guidelines

- You need to consider what you want to annotate and then to define it appropriately

- With multiple annotators it's essential to have a clear set of guidelines for them to follow

- Consistency of annotation is really important for a proper evaluation

# Annotation Guidelines

- People need clear definition of what to annotate in the documents, with examples

- Typically written as a guidelines document

- Piloted first with few annotators, improved, then "real" annotation starts, when all annotators are trained

- Annotation tools may require the definition of a formal DTD (e.g. XML schema)

  - What annotation types are allowed

  - What are their attributes/features and their values

  - Optional vs obligatory; default values

# Annotation Editor

# Annotation Recap

- Adding annotation sets

- Adding annotations

- Resizing them (changing boundaries)

- Deleting

- Changing highlighting colour

- Setting features and their values

- Using the co-reference editor

# Evaluation



"We didn't underperform. You overexpected."

# Performance Evaluation

2 main requirements:

- **Evaluation metric**: mathematically defines how to measure the system's performance against human-annotated gold standard

- **Scoring program**: implements the metric and provides performance measures
  - For each document and over the entire corpus
  - For each type of annotation

# AnnotationDiff

- Graphical comparison of 2 sets of annotations
- Visual diff representation, like tkdiff
- Compares one document at a time, one annotation type at a time

# Annotations are like squirrels…



Annotation Diff helps with "spot the difference"

# Annotation Diff Exercise

- Open the document in-whitbread-10-aug-2001.xml (after processing with ANNIE), create a new Key annotation set and add some new Person annotations there. Add some incorrect annotations as well as correct ones.

- Open the AnnotationDiff: Tools → Annotation Diff or click the icon (green and red pencils)

- Select the name of the document you annotated

- Key contains the manual annotations (select **Key** annotation set)

- Response contains annotations from ANNIE (select **Default** annotation set)

- Select the **Person** annotation

- Click on "Compare"

# Annotation Diff

# A Word about Terminology

- Different communities use different terms when talking about evaluation, because the tasks are a bit different.

- The IE community usually talks about "correct", "spurious" and "missing"

- The IR community usually talks about "true positives", "false positives" and "negatives". They also talk about "false negatives", but you can ignore those.

- Some terminologies assume that one set of annotations is correct ("gold standard")

- Other terminologies do not assume one annotation set is correct

- When measuring inter-annotator agreement, there is no reason to assume one annotator is more correct than the other

# Terminology Comparison

| Gold Standard (IE) | Gold Standard (IR) | Inter-annotator Agreement |
|---|---|---|
| Correct | True Positive | Match |
| Missing | False Negative | Only A (or B) |
| Spurious | False Positive | Only B (or A) |
| Partially Correct | | Overlap |
| | True Negative | |

# **Measuring success**

- In IE, we classify the annotations produced in one of 4 ways:

- **Correct** = things annotated correctly

e.g. annotating "Hamish Cunningham" as a Person

- **Missing** = things not annotated that should have been

e.g. not annotating "Sheffield" as a Location

- **Spurious** = things annotated wrongly

e.g. annotating "Hamish Cunningham" as a Location

- **Partially correct** = the annotation type is correct, but the span is wrong

e,g, annotating just "Cunningham" as a Person (too short) or annotating "Unfortunately Hamish Cunningham" as a Person (too long)

# Finding Precision, Recall and F-measure



scores displayed

# Precision

- How many of the entities your application found were correct?

- Sometimes precision is called **accuracy**

$$Precision = \frac{Correct}{Correct + Spurious}$$

# Recall

- How many of the entities that exist did your application find?

- Sometimes recall is called **coverage**

$$Recall = \frac{Correct}{Correct + Missing}$$

# F-Measure

- Precision and recall tend to trade off against one another

    - If you specify your rules precisely to improve precision, you may get a lower recall

- If you make your rules very general, you get good recall, but low precision

- This makes it difficult to compare applications, or to check whether a change has improved or worsened the results overall

- F-measure combines precision and recall into one measure

# F-Measure

- Also known as the "harmonic mean"

- Usually, precision and recall are equally weighted

- This is known as F1

- To use F1, set the value of the F-measure weight to 1

- This is the default setting

$$F = 2 \cdot \left( \frac{precision \cdot recall}{precision + recall} \right)$$

# Annotation Diff defaults to F1



F-measure weight set to 1

# Statistics can mean what you want them to....

- How we want to measure partially correct annotations may differ, depending on our goal

- In GATE, there are 3 different ways to measure them

- The most usual way is to consider them to be "half right"

- Average: Strict and lenient scores are averaged (this is the same as counting a half weight for every partially correct annotation)

- Strict: Only perfectly matching annotations are counted as correct

- Lenient: Partially matching annotations are counted as correct. This makes your scores look better!

# **Statistics can mean what you want them to....**

Huff (1954) is still the classic text on the subject!

# Strict, Lenient and Average

# Comparing the individual annotations

- In the AnnotationDiff, colour codes indicate whether the annotation pair shown are correct, partially correct, missing (false negative) or spurious (false positive)

- You can sort the columns however you like

# Comparing the annotations



Key annotations

Response annotations

# Corpus Quality Assurance

- Corpus Quality Assurance tool extends the Annotation Diff functionality to the entire corpus, rather than on a single document at a time

- It produces statistics both for the corpus as a whole (Corpus statistics tab) and for each document separately (Document statistics tab)

- It compares two annotation sets, but makes no assumptions about which (if either) set is the gold standard. It just labels them A and B.

- This is because it can be used to measure Inter Annotator Agreement (IAA) where there is no concept of "correct" set

# Try out Corpus Quality Assurance

- Open your hands-on corpus.

- Nine of the documents (ft-*.xml) were manually annotated for you (in the Key AS).

- You annotated one of them (in-whitbread-10-aug-2001.xml).

- Click the Corpus Quality Assurance tab at the bottom of the Display pane.

# Try out Corpus Quality Assurance



- Click the Corpus Quality Assurance tab at the bottom of the Display pane.

# Select Annotation Sets



- Select the annotation sets you wish to compare.

- Click on an annotation set – this will label it set A.

- Now click on another annotation set - this will label it set B

# Select Type

- Select the annotation type to compare

- Select the features to include (if any)

- You can select as many as you want.

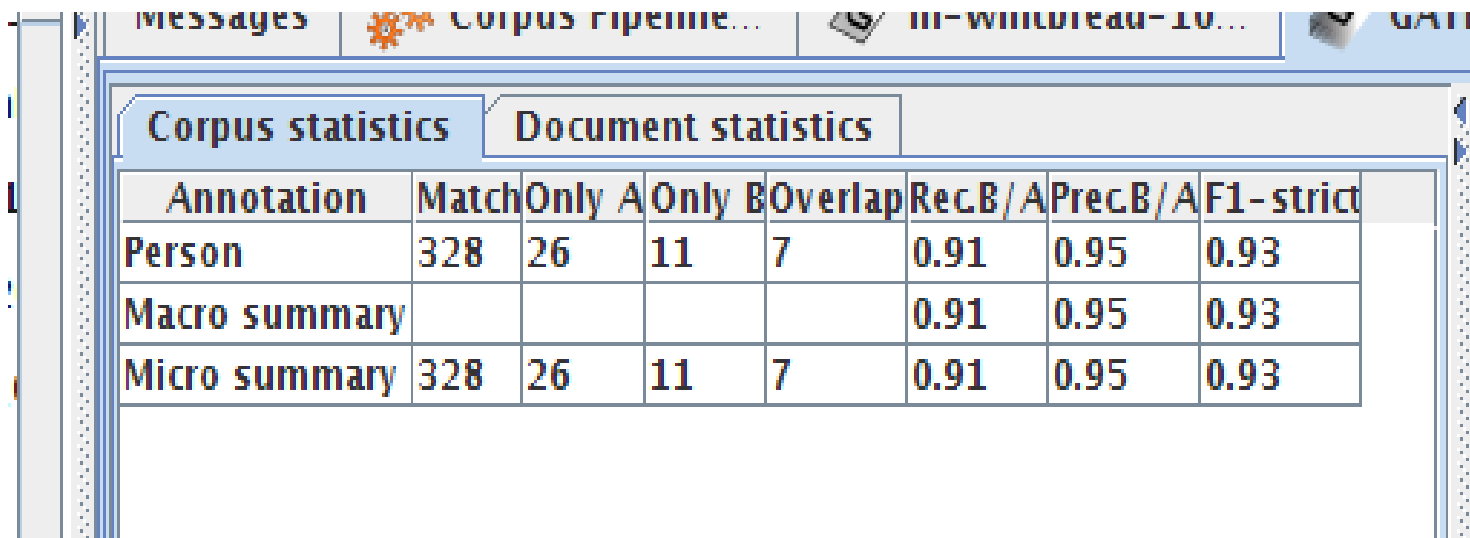- In the "Measures" box, select the kind of F score you want "Strict, Lenient, Average" or any combination of them.

- Select Compare

# Corpus Statistics Tab

| Corpus statistics | Document statistics | | | | | | |
|---|---|---|---|---|---|---|---|
| Annotation | Match | Only A | Only B | Overlap | Rec.B/A | Prec.B/A | F1-strict |
| Person | 328 | 26 | 11 | 7 | 0.91 | 0.95 | 0.93 |
| Macro summary | | | | | 0.91 | 0.95 | 0.93 |
| Micro summary | 328 | 26 | 11 | 7 | 0.91 | 0.95 | 0.93 |

- Each annotation type is listed separately

- Precision, recall and F measure are given for each

- Two summary rows provide micro and macro averages

# Document Statistics Tab

| Document | Match | Only A | Only B | Overlap | Rec.B/A | Prec.B/A | F1-strict |
|----------|-------|--------|--------|---------|---------|----------|-----------|
| in-reed-10-aug-2001.xml_00072 | 10 | 1 | 0 | 0 | 0.91 | 1.00 | 0.95 |
| in-rover-10-aug-2001.xml_00073 | 3 | 0 | 0 | 0 | 1.00 | 1.00 | 1.00 |
| in-scoot-10-aug-2001.xml_00074 | 1 | 0 | 0 | 0 | 1.00 | 1.00 | 1.00 |
| in-shell-cirywire-03-aug-2001.xml_00075 | 7 | 1 | 0 | 0 | 0.88 | 1.00 | 0.93 |
| in-tesco-citywire-07-aug-2001.xml_00076 | 1 | 0 | 0 | 0 | 1.00 | 1.00 | 1.00 |
| in-whitbread-10-aug-2001.xml_00077 | 1 | 0 | 0 | 0 | 1.00 | 1.00 | 1.00 |
| Macro summary | | | | | 0.95 | 0.95 | 0.94 |
| Micro summary | 328 | 26 | 11 | 7 | 0.91 | 0.95 | 0.93 |

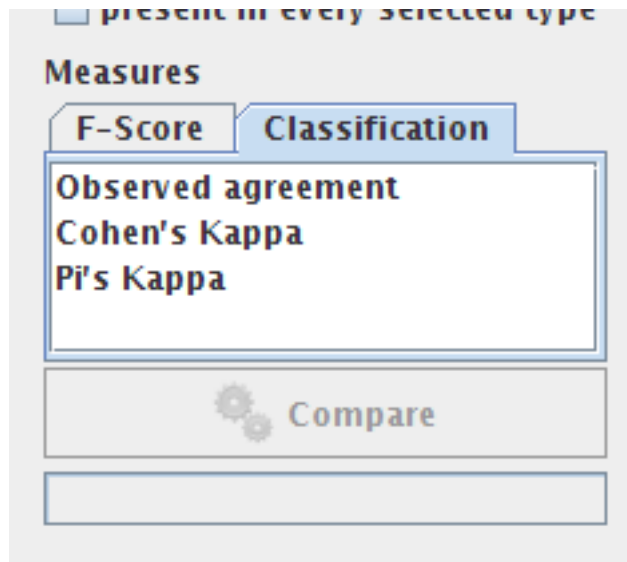Corpus editor  Initialisation Parameters  Corpus Quality Assurance

- Each document is listed separately

- Precision, recall and F measure are given for each

- Two summary rows provide micro and macro averages

# Micro and Macro Averaging

- Micro averaging treats the entire corpus as one big document, for the purposes of calculating precision, recall and F

- Macro averaging takes the average of the rows

# Classification Measures

present in every selected type

**Measures**

| F-Score | Classification |

Observed agreement
Cohen's Kappa
Pi's Kappa

Compare

- By default, Corpus Quality Assurance presents the F-measures

- However, classification measures are also available

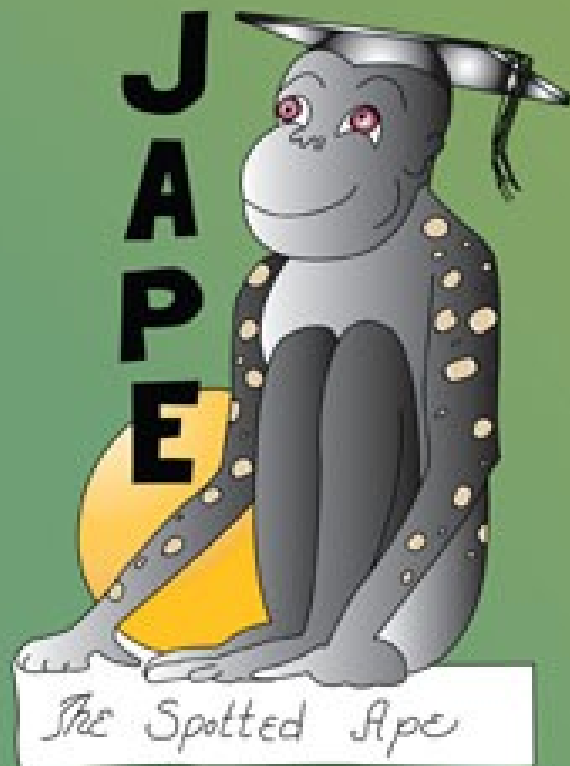- These are not suitable for entity extraction tasks

# Summary

- Module 2 has been devoted to IE and ANNIE

- You should now have a basic understanding of:

    - what IE is

    - how to load and run ANNIE

    - what each of the ANNIE components do

    - how to modify ANNIE components

    - multilingual capabilities of GATE

    - Evaluation

# Tomorrow: introducing JAPE



JAPE, a happy little ape, was always kind and thoughtful. His fine, bright mind helped him find his place in life with an unusual solution to his problem....

# Further exercise: Sentence Splitter variants

- Organisations do not span sentence boundaries, according to the rules used to create them.

- Load the default ANNIE and run it on the document in the directory module2-hands-on/universities

- Look at the Organisation annotations

- Now remove the sentence splitter and replace it with the alternate sentence splitter (see slide on Sentence Splitting variants for details)

- Run ANNIE again and look at the Organisation annotations.

- Can you see the difference?

- Can you understand why? If not, have a look at the relevant Setence annotations.