



Manual Annotation: What is it ? How to do it (properly)?

Karën Fort

karen.fort@loria.fr / <https://members.loria.fr/KFort/>

École thématique d'été "Annotations" - May 31st, 2022

Qual Program

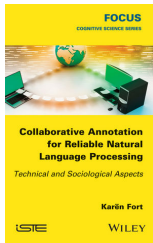
<https://members.loria.fr/KFort/publications/tutorials-summer-schools-etc/>

- ▶ Qual1: now
- ▶ Qual2 (Thursday, 11 am): inter-annotator agreement
- ▶ Qual3 (Friday, 9 am): crowdsourcing

Where I'm talking from

<https://members.loria.fr/KFort/>

► Manual annotation for NLP



► Ethics and NLP



Manuel Annotation and NLP

What is Annotation?

How to do this properly?

Analysing the complexity of an annotation campaign

To finish

Manuel Annotation and NLP

Manuel Annotation in NLP

A notoriously high cost

About language resources longevity

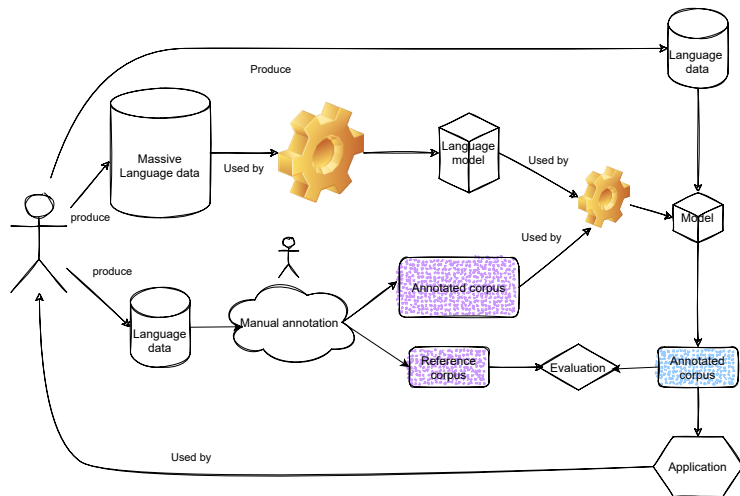
What is Annotation?

How to do this properly?

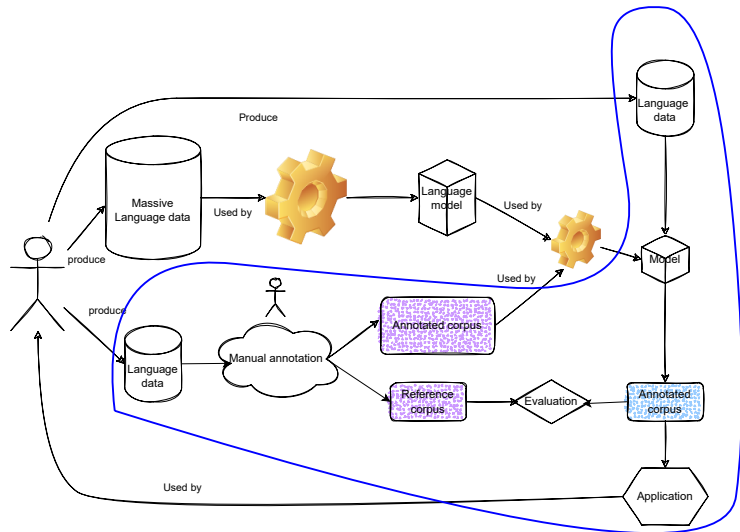
Analysing the complexity of an annotation campaign

To finish

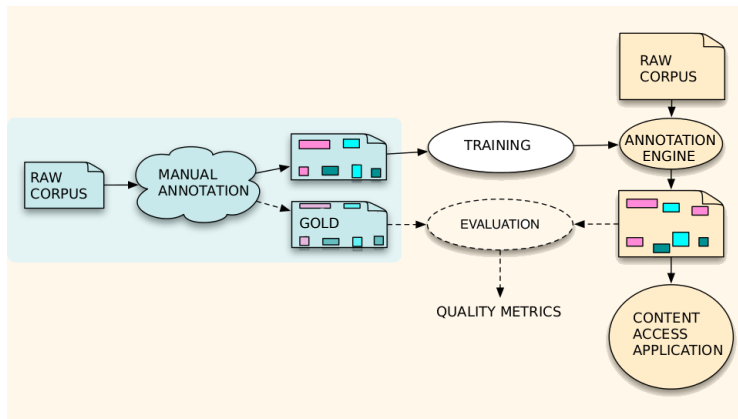
NLP today



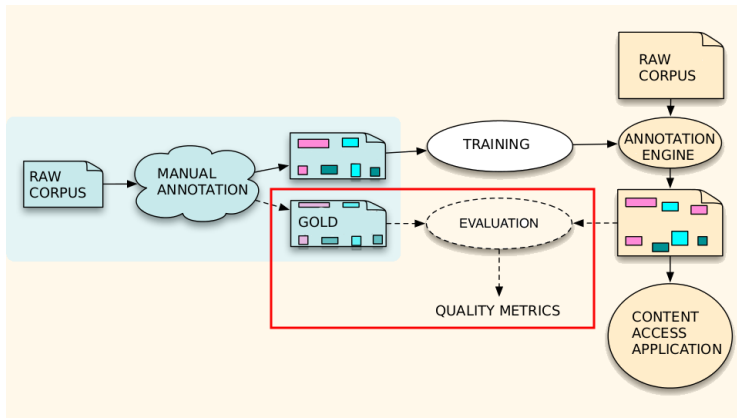
NLP today



Manual annotation in|for NLP



Manual annotation in|for NLP



Why it's important



Ben Hamner  @benhamner · Oct 9



Programming: 10% writing code. 90% figuring out why it doesn't work

Analyzing data and ML: 1% writing code. 9% figuring out why code doesn't work. 90% figuring out what's wrong with the data

 89

 1.9K

 8.7K



Penn Treebank (PTB) [Marcus et al., 1993]

PTB 1:

- ▶ POS-tagging correction: ? units per hour, ? hours a day
- ▶ constituents (syntax) correction: ? units per hour, ? hours a day

Penn Treebank (PTB) [Marcus et al., 1993]

PTB 1:

- ▶ POS-tagging correction: 3,000 units per hour, ? hours a day
- ▶ constituents (syntax) correction: ? units per hour, ? hours a day

Penn Treebank (PTB) [Marcus et al., 1993]

PTB 1:

- ▶ POS-tagging correction: 3,000 units per hour, 3 hours a day
- ▶ constituents (syntax) correction: ? units per hour, ? hours a day

Penn Treebank (PTB) [Marcus et al., 1993]

PTB 1:

- ▶ POS-tagging correction: 3,000 units per hour, 3 hours a day
- ▶ constituents (syntax) correction: 750 units per hour, ? hours a day

Penn Treebank (PTB) [Marcus et al., 1993]

PTB 1:

- ▶ POS-tagging correction: 3,000 units per hour, 3 hours a day
- ▶ constituents (syntax) correction: 750 units per hour, 3 hours a day

Penn Treebank (PTB) [Marcus et al., 1993]

PTB 1:

- ▶ POS-tagging correction: 3,000 units per hour, 3 hours a day
- ▶ constituents (syntax) correction: 750 units per hour, 3 hours a day
- ▶ + learning curve 1 month (POS-tagging) to 2 months (syntax)!

Prague Dependency Treebank [Böhmová et al., 2001]

- ▶ 1996-2004 [Böhmová et al., 2001],
- ▶ built from the CNC (Czech National Corpus),
- ▶ 3 structural levels:
 1. morphological (semi-automatic): 1.8 million tokens
 2. analytical (dependency syntax, with an ad-hoc tool)
 3. tectogrammatical (semantic): 1 million tokens

Prague Dependency Treebank [Böhmová et al., 2001]

Version 1.0:

- ▶ manual annotation of the morphological and analytical levels
- ▶ time: ?
- ▶ nb of persons: ?
- ▶ cost estimate: ?

Prague Dependency Treebank [Böhmová et al., 2001]

Version 1.0:

- ▶ manual annotation of the morphological and analytical levels
- ▶ time: 5 years
- ▶ nb of persons: ?
- ▶ cost estimate: ?

Prague Dependency Treebank [Böhmová et al., 2001]

Version 1.0:

- ▶ manual annotation of the morphological and analytical levels
- ▶ time: 5 years
- ▶ nb of persons: 22 persons, incl. 17 simultaneously during the most demanding periods
- ▶ cost estimate: ?

Prague Dependency Treebank [Böhmová et al., 2001]

Version 1.0:

- ▶ manual annotation of the morphological and analytical levels
- ▶ time: 5 years
- ▶ nb of persons: 22 persons, incl. 17 simultaneously during the most demanding periods
- ▶ cost estimate: \$600,000

GENIA [Kim et al., 2008]

GENIA: 400,000 words annotated in microbiology.

GENIA [Kim et al., 2008]

GENIA: 400,000 words annotated in microbiology.

⇒ 5 part-time annotators, 1 senior coordinator, 1 junior coordinator during 1.5 year [Kim et al., 2008]

ESTER

- ▶ 100 h of transcribed speech (ESTER evaluation campaign, transcription systems, 2008)
- ▶ 1 h of speech = ?

ESTER

- ▶ 100 h of transcribed speech (ESTER evaluation campaign, transcription systems, 2008)
- ▶ 1 h of speech = between 20 and 60 h of transcription work

ESTER

- ▶ 100 h of transcribed speech (ESTER evaluation campaign, transcription systems, 2008)
- ▶ 1 h of speech = between 20 and 60 h of transcription work

⇒ quality has to be high!

Lifespan of annotated corpora

Penn Treebank [Marcus et al., 1993]:

- ▶ created at the beg. of the 90s
- ▶ still used in 2022 (ACL 2022)

vs PARTS tagger [Church, 1988], which is no more used

→ fast evolution of the tools

⇒ annotation should **not** depend on them

Manuel Annotation and NLP

What is Annotation?

Practice

DefinitionS

What for?

How to do this properly?

Analysing the complexity of an annotation campaign

To finish

Practice

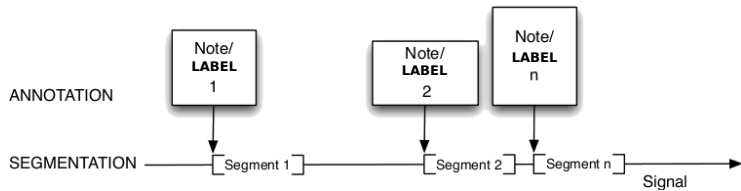
Transcribe what you hear, using any text editor or paper.

Definition

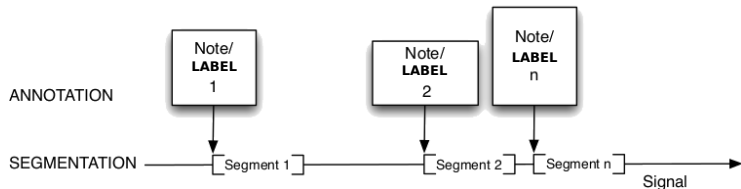
*“[corpus annotation] can be defined as the practice of adding **interpretative**, linguistic information to an electronic corpus of spoken and/or written language data. ‘Annotation’ can also refer to the end-product of this process” [Leech, 1997]*

*“‘Linguistic annotation’ covers any **descriptive** or **analytic** notations applied to raw language data. The basic data may be in the form of time functions - audio, video and/or physiological recordings - or it may be textual.” [Bird and Liberman, 2001]*

Annotation



Annotation



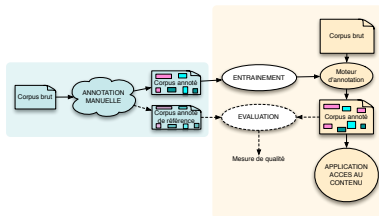
Adding **interpretative** information [Leech, 1997, Habert, 2005]

The application: horizon of the annotation

An annotation is always *task-oriented* [Habert, 2000].

- ▶ direct applicative purpose (summaries of football matches for the football campaign)
- ▶ intermediary application or internal to NLP application (POS-tagging)

[T]he annotations are more useful, the more they are designed to be specific to a particular application [Leech, 2005].



Exercise: annotate soccer match comments

players, teams, actions (goals), relations (passes), etc.

With a huge surprise from the side of Bayern Munich as Van Bommel, the captain, has been removed. He is not even on the substitutes list.

Exercise: annotate soccer match comments

players, teams, actions (goals), relations (passes), etc.

*With a huge surprise from the side of Bayern Munich as Van Bommel, the captain, has been **removed**. He is not even on the substitutes list.*

What is the task, the application aimed at?

summary of match

Van Bommel?

should **not** be annotated

The consensus, at the heart of annotation

One needs to "agree to be able to measure" [Desrosières, 2008]

Annotation is related to **quantification**

Measuring vs quantifying [Desrosières, 2008] :

- ▶ **measuring**: implies a measurable form (eg. the height of Mont Blanc)
- ▶ **quantifying**: implies preliminary conventions of equivalence

The consensus should be equipped:

- ▶ annotation guidelines (12p. for soccer)
- ▶ meetings with the annotators and the campaign manager
- ▶ **evaluate** the consensus (consistency)

Manuel Annotation and NLP

What is Annotation?

How to do this properly?

Good Practises

Theorizing

Analysing the complexity of an annotation campaign

To finish

Leech's 7 maxims [Leech, 1993]

1. It should always be possible to come **back** to initial data (example BC). Note: can be hard after normalization (“l'arbre” → “le arbre”, etc.)

Leech's 7 maxims [Leech, 1993]

1. It should always be possible to come **back** to initial data (example BC). Note: can be hard after normalization (“l'arbre” → “le arbre”, etc.)
2. Annotations should be **extractable** from the text

Leech's 7 maxims [Leech, 1993]

1. It should always be possible to come **back** to initial data (example BC). Note: can be hard after normalization (“l’arbre” → “le arbre”, etc.)
2. Annotations should be **extractable** from the text
3. The annotation procedure should be **documented** (ex: **Brown Corpus annotation guide**, **Penn Tree Bank annotation guide**)

Leech's 7 maxims [Leech, 1993]

1. It should always be possible to come **back** to initial data (example BC). Note: can be hard after normalization (“l'arbre” → “le arbre”, etc.)
2. Annotations should be **extractable** from the text
3. The annotation procedure should be **documented** (ex: **Brown Corpus annotation guide**, **Penn Tree Bank annotation guide**)
4. Mention should be made of the **annotator(s)** and the way annotation was made (manual/automatic annotation, number of annotators, manually corrected/uncorrected...)

Leech's 7 maxims [Leech, 1993]

1. It should always be possible to come **back** to initial data (example BC). Note: can be hard after normalization (“l’arbre” → “le arbre”, etc.)
2. Annotations should be **extractable** from the text
3. The annotation procedure should be **documented** (ex: **Brown Corpus annotation guide**, **Penn Tree Bank annotation guide**)
4. Mention should be made of the **annotator(s)** and the way annotation was made (manual/automatic annotation, number of annotators, manually corrected/uncorrected...)
5. Annotation is an act of **interpretation** (cannot be infallible)

Leech's 7 maxims [Leech, 1993]

1. It should always be possible to come **back** to initial data (example BC). Note: can be hard after normalization (“l’arbre” → “le arbre”, etc.)
2. Annotations should be **extractable** from the text
3. The annotation procedure should be **documented** (ex: **Brown Corpus annotation guide**, **Penn Tree Bank annotation guide**)
4. Mention should be made of the **annotator(s)** and the way annotation was made (manual/automatic annotation, number of annotators, manually corrected/uncorrected...)
5. Annotation is an act of **interpretation** (cannot be infallible)
6. Annotation schemas should be as **independent** as possible on formalisms

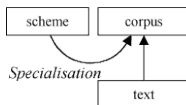
Leech's 7 maxims [Leech, 1993]

1. It should always be possible to come **back** to initial data (example BC). Note: can be hard after normalization (“l’arbre” → “le arbre”, etc.)
2. Annotations should be **extractable** from the text
3. The annotation procedure should be **documented** (ex: **Brown Corpus annotation guide**, **Penn Tree Bank annotation guide**)
4. Mention should be made of the **annotator(s)** and the way annotation was made (manual/automatic annotation, number of annotators, manually corrected/uncorrected...)
5. Annotation is an act of **interpretation** (cannot be infallible)
6. Annotation schemas should be as **independent** as possible on formalisms
7. No annotation schema should consider itself a standard (it possibly becomes one)

Different points of view

“you only get out what you put in” [Wallis, 2007]

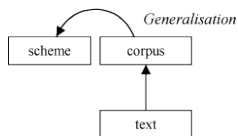
Model-based approach



Knowledge is in the annotation schema \Rightarrow corpus comes after

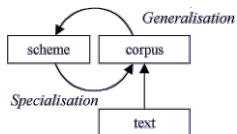
Everything is in the annotation!

Corpus-based approach



Knowledge is in the text \Rightarrow the corpus comes first [Sinclair]

Third way?



The knowledge is in the annotation schema **and** in the corpus

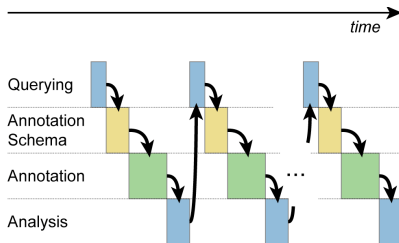
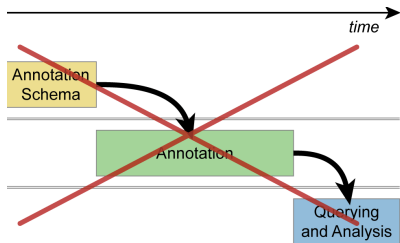
Annotation by cycles

- ▶ new observations generalize hypotheses
 - ▶ theory allows to interpret and classify information

 - ▶ **evolving** cycles: each cycle improves the knowledge by refining and testing the theories on real data
- ⇒ a **more precise** representation of the corpus is built and a **more sophisticated** system is produced

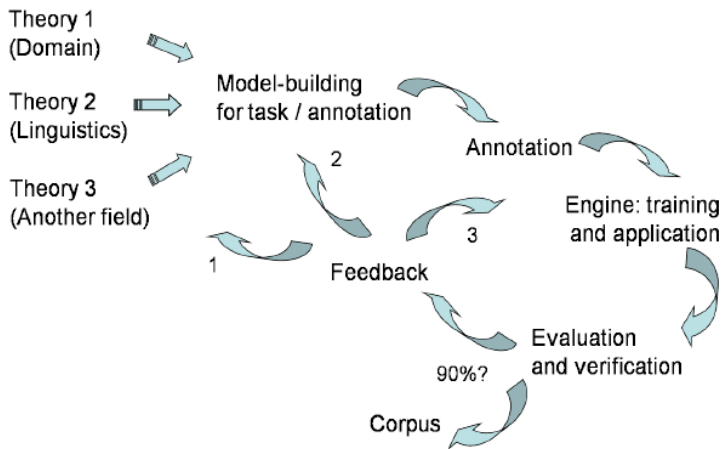
Agile Annotation

integrating evaluation

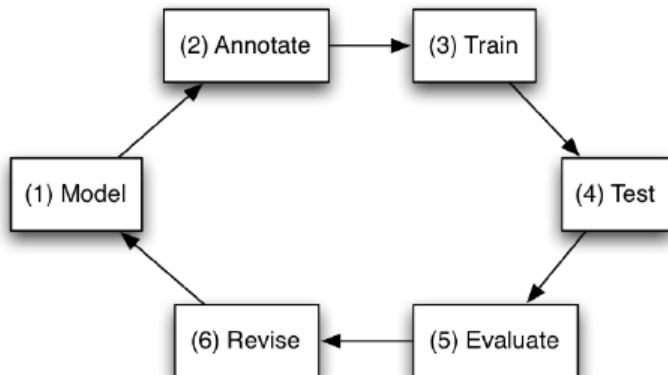


Traditional annotation phases (left) and cycles of agile annotation (right). Reproduction of Figure 2 from [Voormann and Gut, 2008]

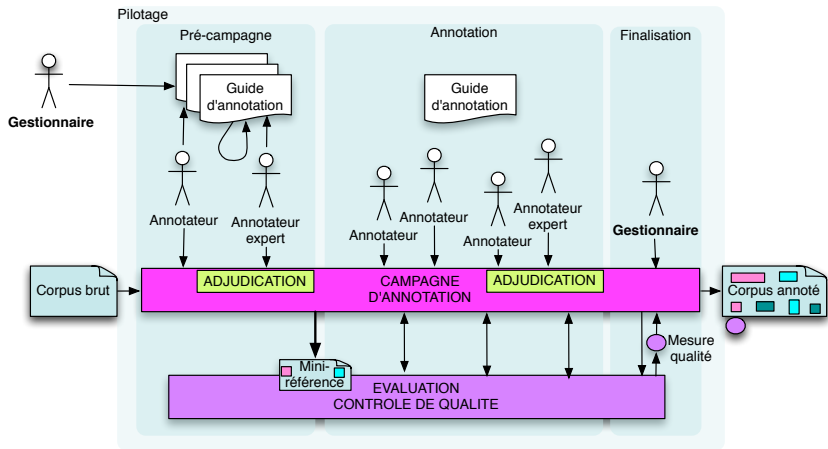
Generic annotation pipeline [Hovy and Lavid, 2010]



MATTER cycle [Pustejovsky and Stubbs, 2012]



Towards "annotation engineering" [Fort, 2012]



Manuel Annotation and NLP

What is Annotation?

How to do this properly?

Analysing the complexity of an annotation campaign

- What to annotate?

- How to annotate?

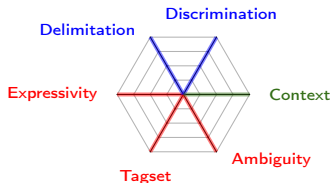
- Weight of the context

- Synthesis

To finish

Complexity dimensions

- ▶ 5 independent dimensions:
 - ▶ 2 related to the localisation of annotations
 - ▶ 3 related to the characterisation of annotations
- ▶ 1 not independent: the context



- ▶ Scale from **0** (null complexity) to **1** (maximal complexity) to allow for the comparison between campaigns
- ▶ Independent from the volume to annotate and the number of annotators

Example: gene renaming

1. Identification of gene names in the source signal:

*The **yppB** gene complemented the defect of the recG40 strain. **yppB** and **yppC** and their respective null alleles were termed “**recU**” and “**recU1**” (recU:cat) and “**recS**” and “**recS1**” (recS:cat), respectively.*

2. Identification of gene couples expressing a renaming relation:

*The **yppB** gene complemented the defect of the recG40 strain. **yppB** and **yppC** and their respective null alleles were termed “**recU**” and “**recU1**” (recU:cat) and “**recS**” and “**recS1**” (recS:cat), respectively.*

Discrimination

Parts-of-speech [Marcus et al., 1993], pre-annotated :

I/PRP do/VBP n't/RB feel/VB very/RB ferocious/JJ ./.

Gene renaming [Fort et al., 2012], no pre-annotation:

The yppB:cat and ypbC:cat null alleles rendered cells sensitive to DNA-damaging agents, impaired plasmid transformation (25- and 100-fold), and moderately affected chromosomal transformation when present in an otherwise Rec+ B. subtilis strain. The yppB gene complemented the defect of the recG40 strain. yppB and ypbC and their respective null alleles were termed recU and “recU1” (recU:cat) and recS and “recS1” (recS:cat), respectively. The recU and recS mutations were introduced into rec-deficient strains representative of the alpha (recF), beta (addA5 addB72), gamma (recH342), and epsilon (recG40) epistatic groups.

Discrimination

Parts-of-speech [Marcus et al., 1993], pre-annotated :

I/PRP do/VBP n't/RB feel/VB very/RB ferocious/JJ ./.

Gene renaming [Fort et al., 2012], no pre-annotation:

The yppB:cat and ypbC:cat null alleles rendered cells sensitive to DNA-damaging agents, impaired plasmid transformation (25- and 100-fold), and moderately affected chromosomal transformation when present in an otherwise Rec+ B. subtilis strain. The yppB gene complemented the defect of the recG40 strain. yppB and ypbC and their respective null alleles were termed recU and "recU1" (recU:cat) and recS and "recS1" (recS:cat), respectively. The recU and recS mutations were introduced into rec-deficient strains representative of the alpha (recF), beta (addA5 addB72), gamma (recH342), and epsilon (recG40) epistatic groups.

⇒ **more difficult** if the units to annotate are scattered, in particular if the segmentation is not obvious.

Discrimination

The discrimination weight is all the more high as the proportion of what *should* be annotated as compared to what *could* be annotated is low.

Definition

$$Discrimination(Flow) = 1 - \frac{|Annotations(Flow)|}{\sum_{i=1}^{LevelSeg} |UnitsObtainedBySeg_i(Flow)|}$$

⇒ Need for a [reference segmentation](#)

Parts-of-speech[Marcus et al., 1993] :

I/PRP do/VBP n't/RB feel/VB very/RB ferocious/JJ ./.

$$Discrimination_{PTB_{POS}} = 0$$

Gene renaming[Fort et al., 2012] :

The yppB:cat and ypbC:cat null alleles rendered cells sensitive to DNA-damaging agents, impaired plasmid transformation (25- and 100-fold), and moderately affected chromosomal transformation when present in an otherwise Rec+ B. subtilis strain. The yppB gene complemented the defect of the recG40 strain. yppB and ypbC and their respective null alleles were termed recU and “recU1” (recU:cat) and recS and “recS1” (recS:cat), respectively. The recU and recS mutations were introduced into rec-deficient strains representative of the alpha (recF), beta (addA5 addB72), gamma (recH342), and epsilon (recG40) epistatic groups.

$$Discrimination_{Identification} = 0,9$$

$$Discrimination_{Renaming} = 0,95$$

Boundaries delimitation

- ▶ **extending** or **shrinking** the discriminated unit:
Madame Chirac → *Monsieur et Madame Chirac*

Boundaries delimitation

- ▶ **extending** or **shrinking** the discriminated unit:
Madame Chirac → *Monsieur et Madame Chirac*
- ▶ **decompose** a discriminated unit into several elements:
le préfet Érignac → *le préfet Érignac*

Boundaries delimitation

- ▶ **extending** or **shrinking** the discriminated unit:
Madame Chirac → *Monsieur et Madame Chirac*
- ▶ **decompose** a discriminated unit into several elements:
le préfet Érignac → *le **préfet** **Érignac***
- ▶ or **group** together several discriminated units into one unique annotation:
Sa Majesté
le roi Mohamed VI → ***Sa Majesté le roi Mohamed VI***

Boundaries delimitation

Definition

$$\text{Delimitation}(\text{Flow}) = \min \left(\frac{\text{Substitutions} + \text{Additions} + \text{Deletions}}{|\text{Annotations}(\text{Flow})|}, 1 \right)$$

$$\text{Delimitation}_{\text{Identification}} = 0$$

$$\text{Delimitation}_{\text{Renaming}} = 0$$

$$\text{Delimitation}_{\text{PTB}_{\text{POS}}} = 0$$

$$\text{Délimitation}_{\text{EN}_{\text{TypesSubtypes}}} = 1$$

$$\text{Délimitation}_{\text{EN}_{\text{Components}}} = 0,3$$

Expressiveness of the annotation language

Definition

The degrees of expressiveness of the annotation language are the following:

- ▶ 0.25: type languages
- ▶ 0.5: relational languages of arity 2
- ▶ 0.75: relational languages of arity higher than 2
- ▶ 1: higher-order languages

$$\text{Expressiveness}_{\text{Identification}} = 0.25$$

$$\text{Expressiveness}_{\text{Renaming}} = 0.25$$

Dimension of the tagset

| Person | | | Function | | |
|---|---|--|---|---|--|
| <i>pers.ind</i> (individual person) | | <i>pers.coll</i> (group of persons) | <i>func.ind</i> (individual function) | | <i>func.coll</i> (collectivity of functions) |
| Location | | | Production | | |
| administrative (<i>loc.adm.town</i> , <i>loc.adm.reg</i> , <i>loc.adm.nat</i> , <i>loc.adm.sup</i>) | physical (<i>loc.phys.geo</i> , <i>loc.phys.hydro</i> , <i>loc.phys.astro</i>) | facilities (<i>loc.fac</i>), oronyms (<i>loc.oro</i>), address (<i>loc.add.phys</i> , <i>loc.add.elec</i>) | <i>prod.object</i> (manufactured object) | <i>prod.serv</i> (transportation route) | <i>prod.fin</i> (financial products) |
| | | | <i>prod.doctr</i> (doctrine) | <i>prod.rule</i> (law) | <i>prod.soft</i> (software) |
| | | | <i>prod.art</i> | <i>prod.media</i> | <i>prod.award</i> |
| Organization | | | Time | | |
| <i>org.adm</i> (administration) | | <i>org.ent</i> (services) | <i>time.date.abs</i> (absolute date), <i>time.date.rel</i> (relative date) | <i>time.hour.abs</i> (absolute hour), <i>time.hour.rel</i> (relative hour) | |
| Amount | | | | | |
| <i>amount</i> (with unit or general object), including duration | | | | | |

Types and sub-types used for structured NE annotation [Grouin et al., 2011]

Dimension of the tagset

| | | | | | |
|--|--|--|---|---|--------------------------------------|
| Person | | | Function | | |
| <i>pers.ind</i> (individual person) | <i>pers.coll</i> (group of persons) | | <i>func.ind</i> (individual function) | <i>func.coll</i> (collectivity of functions) | |
| Location | | | Production | | |
| <i>administrative</i> (<i>loc.adm.town</i> , <i>loc.adm.reg</i> , <i>loc.adm.nat</i> , <i>loc.adm.sup</i>) | physical (<i>loc.phys.geo</i> , <i>loc.phys.hydro</i> , <i>loc.phys.astro</i>) | facilities (<i>loc.fac</i>), oronyms (<i>loc.oro</i>), address (<i>loc.add.phys</i> , <i>loc.add.elec</i>) | <i>prod.object</i> (manufactured object) | <i>prod.serv</i> (transportation route) | <i>prod.fin</i> (financial products) |
| | | | <i>prod.doctr</i> (doctrine) | <i>prod.rule</i> (law) | <i>prod.soft</i> (software) |
| | | | <i>prod.art</i> | <i>prod.media</i> | <i>prod.award</i> |
| Organization | | | Time | | |
| <i>org.adm</i> (administration) | <i>org.ent</i> (services) | | <i>time.date.abs</i> (absolute date), <i>time.date.rel</i> (relative date) | <i>time.hour.abs</i> (absolute hour), <i>time.hour.rel</i> (relative hour) | |
| Amount | | | | | |
| <i>amount</i> (with unit or general object), including duration | | | | | |

Level 1: *pers*, *func*, *loc*, *prod*, *org*, *time*, *amount* → 7 possibilities (degree of freedom = 6).

Dimension of the tagset

| Person | | | Function | | |
|---|---|--|--|--|---|
| <i>pers.ind</i> (individual person) | <i>pers.coll</i> (group of persons) | | <i>func.ind</i> (individual function) | <i>func.coll</i> (collectivity of functions) | |
| Location | | | Production | | |
| <i>loc.adm.town</i> , <i>loc.adm.reg</i> , <i>loc.adm.nat</i> , <i>loc.adm.sup</i> | physical (<i>loc.phys.geo</i> , <i>loc.phys.hydro</i> , <i>loc.phys.astro</i>) | facilities (<i>loc.fac</i>), oronyms (<i>loc.oro</i>), address (<i>loc.add.phys</i> , <i>loc.add.elec</i>) | <i>prod.object</i> (manufactured object) | <i>prod.serv</i> (transportation route) | <i>prod.fin</i> (financial products) |
| | | | <i>prod.doctr</i> (doctrine) | <i>prod.rule</i> (law) | <i>prod.soft</i> (software) |
| | | | <i>prod.art</i> | <i>prod.media</i> | <i>prod.award</i> |
| Organization | | | Time | | |
| <i>org.adm</i> (administration) | | <i>org.ent</i> (services) | <i>time.date.abs</i> (absolute date), <i>time.date.rel</i> (relative date) | <i>time.hour.abs</i> (absolute hour), <i>time.hour.rel</i> (relative hour) | |
| Amount | | | | | |
| <i>amount</i> (with unit or general object), including duration | | | | | |

Level 1: *pers*, *func*, *loc*, *prod*, *org*, *time*, *amount* → 7 possibilities (degree of freedom = 6).

Level 2: *prod.object*, *prod.serv*, *prod.fin*, *prod.soft*, *prod.doctr*, *prod.rule*, *prod.art*, *prod.media*, *prod.award* → 9 possibilities (degree of freedom = 8).

Dimension of the tagset

| Person | | | Function | | |
|---|---|--|---|---|---|
| <i>pers.ind</i> (individual person) | <i>pers.coll</i> (group of persons) | | <i>func.ind</i> (individual function) | <i>func.coll</i> (collectivity of functions) | |
| Location | | | Production | | |
| administrative (<i>loc.adm.town</i> , <i>loc.adm.reg</i> , <i>loc.adm.nat</i> , <i>loc.adm.sup</i>) | physical (<i>loc.phys.geo</i> , <i>loc.phys.hydro</i> , <i>loc.phys.astro</i>) | facilities (<i>loc.fac</i>), oronyms (<i>loc.oro</i>), address (<i>loc.add.phys</i> , <i>loc.add.elec</i>) | <i>prod.object</i> (manufactured object) | <i>prod.serv</i> (transportation route) | <i>prod.fin</i> (financial products) |
| | | | <i>prod.doctr</i> (doctrine) | <i>prod.rule</i> (law) | <i>prod.soft</i> (software) |
| | | | <i>prod.art</i> | <i>prod.media</i> | <i>prod.award</i> |
| Organization | | | Time | | |
| <i>org.adm</i> (administration) | <i>org.ent</i> (services) | | <i>time.date.abs</i> (absolute date), <i>time.date.rel</i> (relative date) | <i>time.hour.abs</i> (absolute hour), <i>time.hour.rel</i> (relative hour) | |
| Amount | | | | | |
| <i>amount</i> (with unit or general object), including duration | | | | | |

Level 1: *pers*, *func*, *loc*, *prod*, *org*, *time*, *amount* → 7 possibilities (degree of freedom = 6).

Level 2: *prod.object*, *prod.serv*, *prod.fin*, *prod.soft*, *prod.doctr*, *prod.rule*, *prod.art*, *prod.media*, *prod.award* → 9 possibilities (degree of freedom = 8).

Level 3: *loc.adm.town*, *loc.adm.reg*, *loc.adm.nat*, *loc.adm.sup* → 4 possibilities (degree of freedom = 3).

Dimension of the tagset

Degree of freedom

$$\nu = \nu_1 + \nu_2 + \dots + \nu_m$$

where ν_i is the maximal degree of freedom the annotator has when choosing the i^{th} sub-type ($\nu_i = n_i - 1$).

Dimension of the tagset

$$\text{Dimension}(\text{Flow}) = \min\left(\frac{\nu}{\tau}, 1\right)$$

where τ is the threshold from which we consider the tagset to be very large (experimentally determined).

$$\begin{aligned}\text{Dimension}_{\text{Identification}} &= 0 \\ \text{Dimension}_{\text{Renaming}} &= 0.04 \\ \text{Dimension}_{\text{NE}_{\text{TypesSubtypes}}} &= 0.34\end{aligned}$$

Degree of ambiguity: residual ambiguity

Using the traces left by the annotators:



[...] *<EukVirus>3CDproM</EukVirus> can process both structural and nonstructural precursors of the <EukVirus uncertainty-type = "too-generic"><taxon>poliovirus</taxon> polyprotein</EukVirus> [...].*

Définition

$$AmbiguityRes(Flow) = \frac{|Annotations_{amb}|}{|Annotations|}$$

$$AmbiguityRes_{Identification} = 0.04$$

$$AmbiguityRes_{Renaming} = 0.02$$

Degree of ambiguity: theoretical ambiguity

Proportion of the units to annotate that corresponds to ambiguous vocables.

Definition

$$AmbiguityTh(Flow) = \frac{\sum_{voc_i=1}^{|Voc(Flow)|} (Ambig(voc_i) * freq(voc_i, Flow))}{|Units(Flow)|}$$

with

$$Ambig(voc_i) = \begin{cases} 1 & \text{if } |Tags(voc_i)| > 1 \\ 0 & \text{else} \end{cases}$$

$$AmbiguityTh_{Identification} = 0.01$$

→ Does not apply to renaming relations

Context to take into account

- ▶ **size of the window** to take into account in the source signal:

- ▶ The sentence:

I/PRP do/VBP n't/RB feel/VB very/RB ferocious/JJ ./.

- ▶ ... or more:

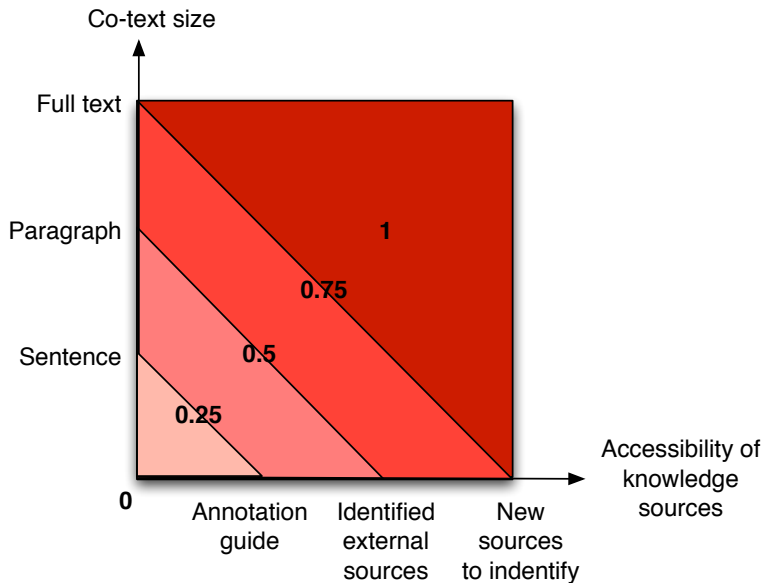
Fabien Lévêque : C'est bien fait , avec **Gouffran** maintenant . **Gouffran** qui va tenter sa chance , et ça fait le but . Le but !

Xavier Gravelaine : Oh la la la la !

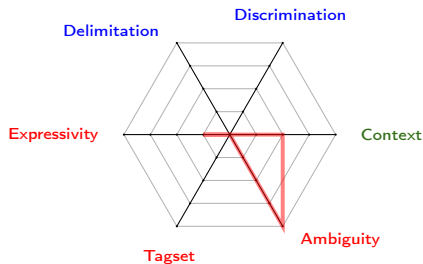
Fabien Lévêque : Et le but du plus breton des **Girondins** , C'est **Yoann Gourcuif** qui vient mettre un quatrième but ici au **stade de France** . Le cauchemar continue pour le **VOC** . Quatre à zéro en faveur des **Girondins** .

- ▶ number of **knowledge elements** to be rallied or degree of accessibility of the knowledge sources that are consulted:
 - ▶ annotation guidelines
 - ▶ nomenclatures (Swiss-Prot)
 - ▶ new sources to be found (Wikipedia, etc.)

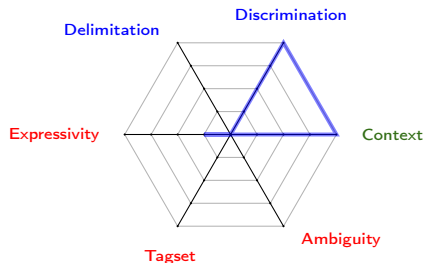
Weight of the context



Synthesis of the complexity dimensions



Classification of it pronouns as anaphoric or impersonal



Gene names identification

What's next: Qual Program

<https://members.loria.fr/KFort/publications/tutorials-summer-schools-etc/>

- ▶ Qual1: now
- ▶ Qual2 (Thursday, 11 am): inter-annotator agreement
- ▶ Qual3 (Friday, 9 am): crowdsourcing

Manuel Annotation and NLP

What is Annotation?

How to do this properly?

Analysing the complexity of an annotation campaign

To finish

WYHTR: What You Have To Remember



Manual annotation and NLP:

- ▶ usage
- ▶ cost

Manual annotation:

- ▶ definition
- ▶ organization
- ▶ complexity grid

A bit of reading

The Manual Annotation Complexity Grid

Modeling the Complexity of Manual Annotation Tasks: a Grid of Analysis

Karën Fort, Adeline Nazarenko, Sophie Rosset

International Conference on Computational Linguistics, Dec 2012, Mumbai, India. pp.895–910

https://hal.archives-ouvertes.fr/hal-00769631/file/coling2012_Complexity_KF_30102012.pdf



Bird, S. and Liberman, M. (2001).

A formal framework for linguistic annotation.

Speech Communication, 33(1-2):23–60.



Böhmová, A., Hajič, J., Hajičová, E., and Hladká, B. (2001).

The prague dependency treebank: Three-level annotation scenario.

In Abeillé, A., editor, Treebanks: Building and Using Syntactically Annotated Corpora. Kluwer Academic Publishers.



Church, K. W. (1988).

A stochastic parts program and noun phrase parser for unrestricted text.

In Proceedings of the Second Conference on Applied Natural Language Processing, ANLC '88, pages 136–143, Stroudsburg, PA, USA. Association for Computational Linguistics.



Desrosières, A. (2008).

Pour une sociologie historique de la quantification : L'Argument statistique.

Presses de l'école des Mines de Paris.



Fort, K. (2012).

Les ressources annotées, un enjeu pour l'analyse de contenu :
vers une méthodologie de l'annotation manuelle de corpus.
PhD thesis, Université Paris XIII, LIPN, INIST-CNRS.



Fort, K., François, C., Galibert, O., and Ghribi, M. (2012).
Analyzing the impact of prevalence on the evaluation of a
manual annotation campaign.

In Proceedings of the International Conference on Language
Resources and Evaluation (LREC), Istanbul, Turquie.
7 pages.



Grouin, C., Rosset, S., Zweigenbaum, P., Fort, K., Galibert, O.,
and Quintard, L. (2011).

Proposal for an extension of traditional named entities: From
guidelines to evaluation, an overview.

In Proceedings of the 5th Linguistic Annotation Workshop,
pages 92–100, Portland, Oregon, USA.
Poster.



Habert, B. (2000).

Corpus. Méthodologie et applications linguistiques, chapter
Détournements d'annotation : armer la main et le regard,
pages 106–120.

Champion and Presses Universitaires de Perpignan.



Habert, B. (2005).

Portrait de linguiste(s) à l'instrument.

Texto!, vol. X(4).



Hovy, E. H. and Lavid, J. M. (2010).

Towards a "science" of corpus annotation: A new
methodological challenge for corpus linguistics.

International Journal of Translation Studies, 22(1).



Kim, J.-D., Ohta, T., and Tsujii, J. (2008).

Corpus annotation for mining biomedical events from literature.

BMC Bioinformatics, 9(1):10.



Leech, G. (1993).

Corpus annotation schemes.

Literary and Linguistic Computing, 8(4):275–281.



Leech, G. (1997).

Corpus annotation: Linguistic information from computer text corpora, chapter Introducing corpus annotation, pages 1–18.

Longman, Londres, Angleterre.



Leech, G. (2005).

Developing Linguistic Corpora: a Guide to Good Practice, chapter Adding Linguistic Annotation, pages 17–29.

Oxford: Oxbow Books.



Marcus, M., Santorini, B., and Marcinkiewicz, M. A. (1993).

Building a large annotated corpus of English : The Penn Treebank.

Computational Linguistics, 19(2):313–330.



Pustejovsky, J. and Stubbs, A. (2012).

Natural Language Annotation for Machine Learning.

O'Reilly.



Voormann, H. and Gut, U. (2008).

Agile corpus creation.

Corpus Linguistics and Linguistic Theory, 4(2):235–251.



Wallis, S. (2007).

Annotating Variation and Change, chapter Annotation,
Retrieval and Experimentation.

Varieng, University of Helsinki, Helsinki, Finland.