

# **Text linguistic foundations of language documentation**

Dafydd Gibbon

Universität Bielefeld

2015-11-02

Abidjan-Bielefeld Cooperation

## The problem to be solved

How is professional documentation created?

What is professional documentation?

- What is its structure?
- What is its physical form?

Meetings:

1. Basics: text linguistics and document creation
2. Documenting spoken language
3. Recording and annotation
4. Language documentation report

# Language documents and their creation

Language documentation is concerned with multimodal documents (artefacts such as texts, audio and video recordings, databases) and their creation with manual and computational methods.

The multimodal documents are designed to have standardised properties pertaining to

- reusability (for purposes other than the original purpose)
- sustainability (for survival over long periods of time)
- interoperability (for use in practical environments)
- ubiquity (for general accessibility)
- coherence (for organising documents and databases)

These properties are supported by using a text linguistic frame of reference.

# What is language documentation?

Language documentation is a manual and computational methodology and a product, with a broader linguistic context:

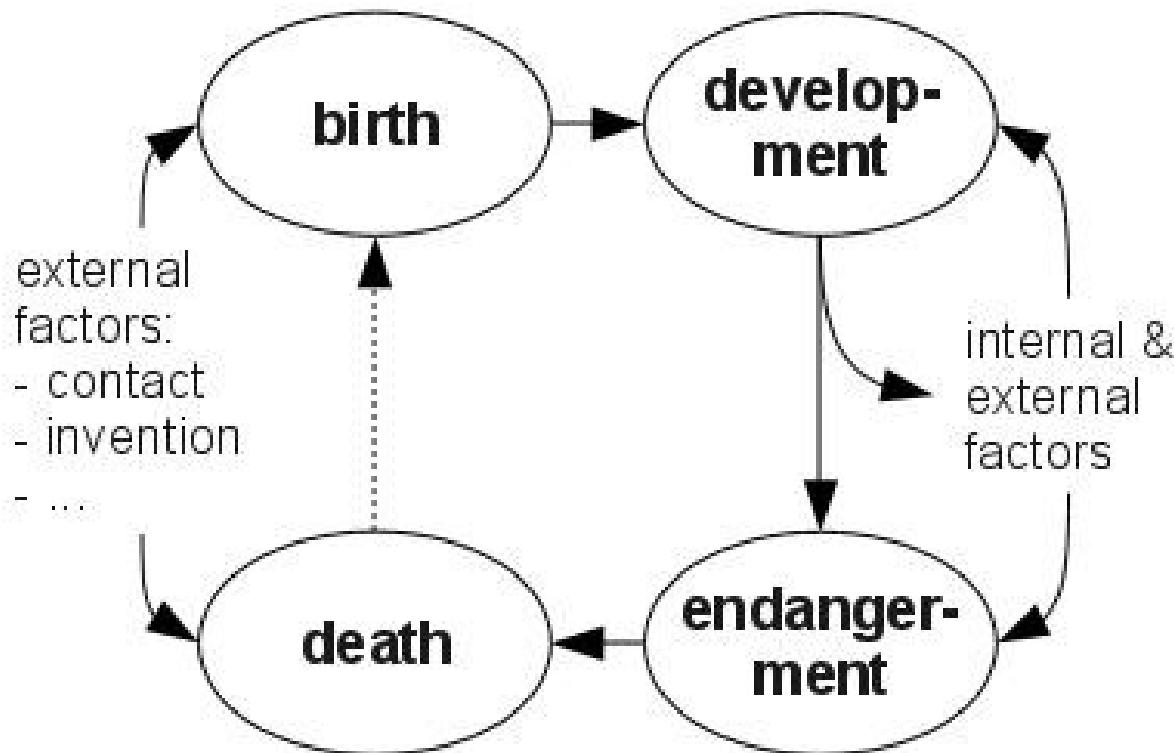
- Discovery methods:
  - documentation: recording of observed data
  - description: generalisation over documentation
  - explanation: predictive theories of descriptions
- Evaluation methods:
  - testing of explanation, description, documentation
  - extension of results to new data
- Application methods (manual and computational):
  - Storage media: databases for
    - standardisation, sustainability, reusability, interoperability
  - Production of objects with tools:
    - books and digital media/texts, dictionaries, grammars

# What is language documentation?

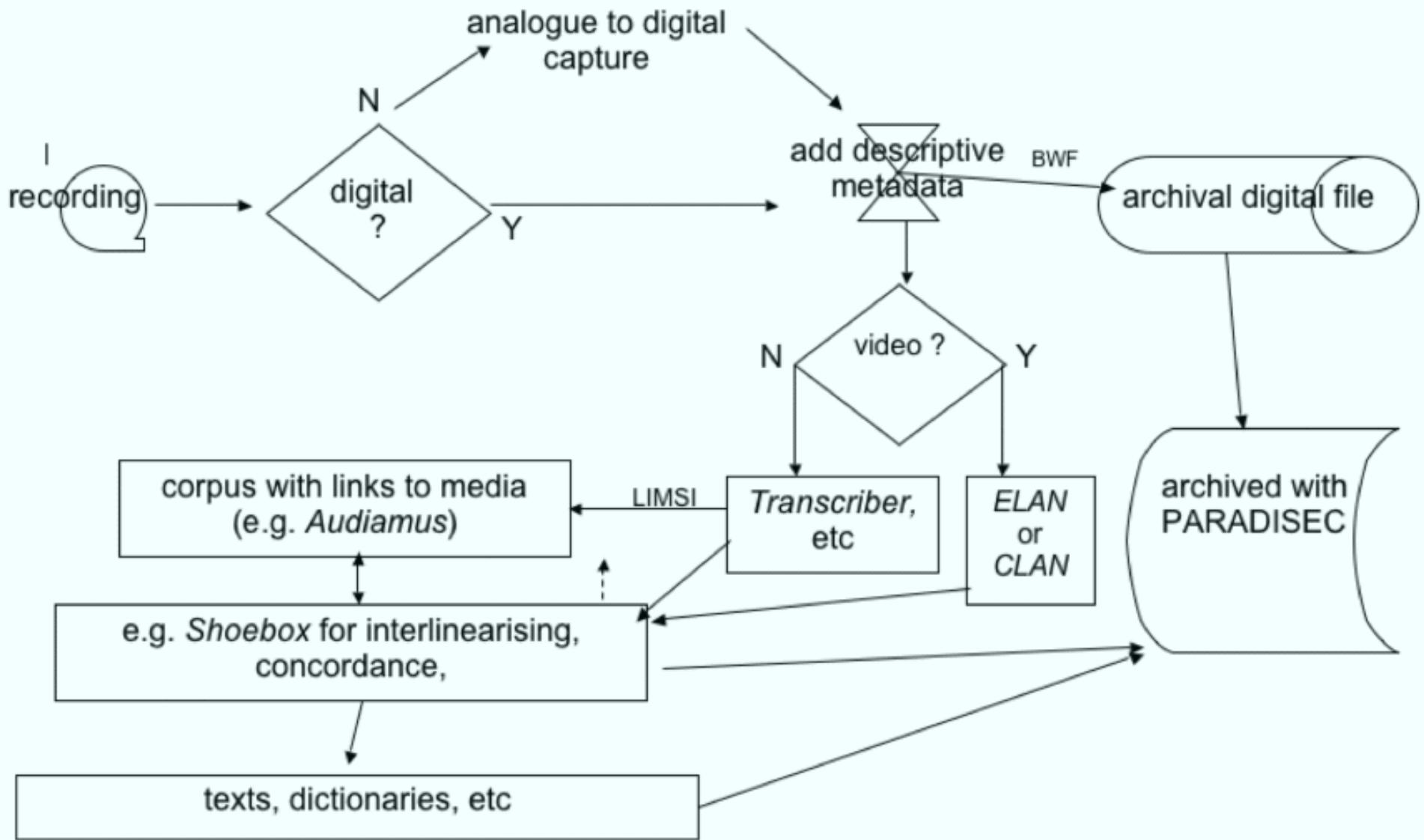
There are many approaches to language documentation:

- field linguistics
- laboratory linguistics
- armchair linguistics
- strongly empiricist
- qualitative and quantitative methods
- observations rather than generalisations
- digital and multimedia products:
  - transcriptions and texts
  - dictionaries
  - grammars
  - audio: recordings of many different data types
  - video: photos, diagrammes, maps, videos

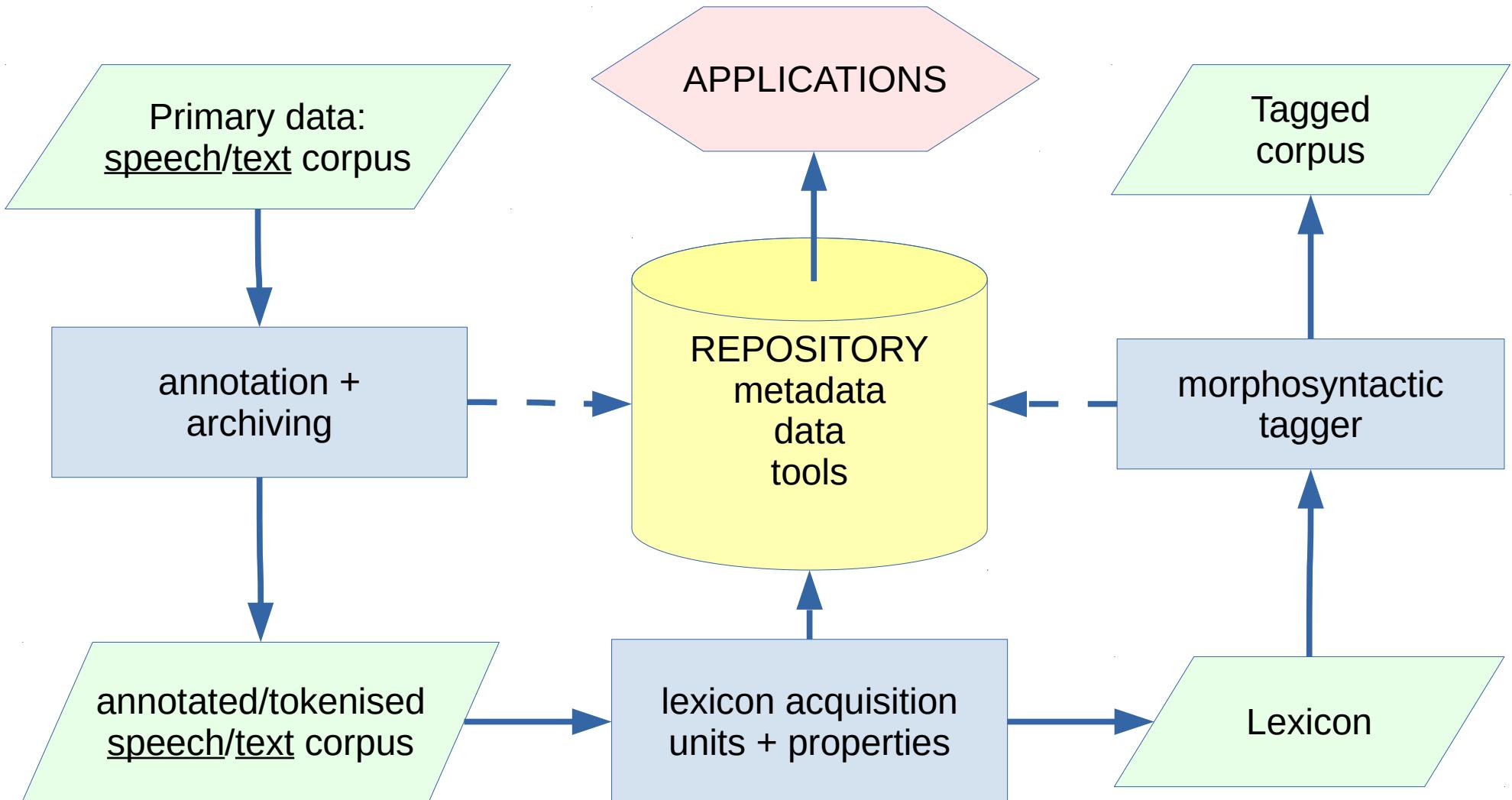
# Documenting the language life-cycle



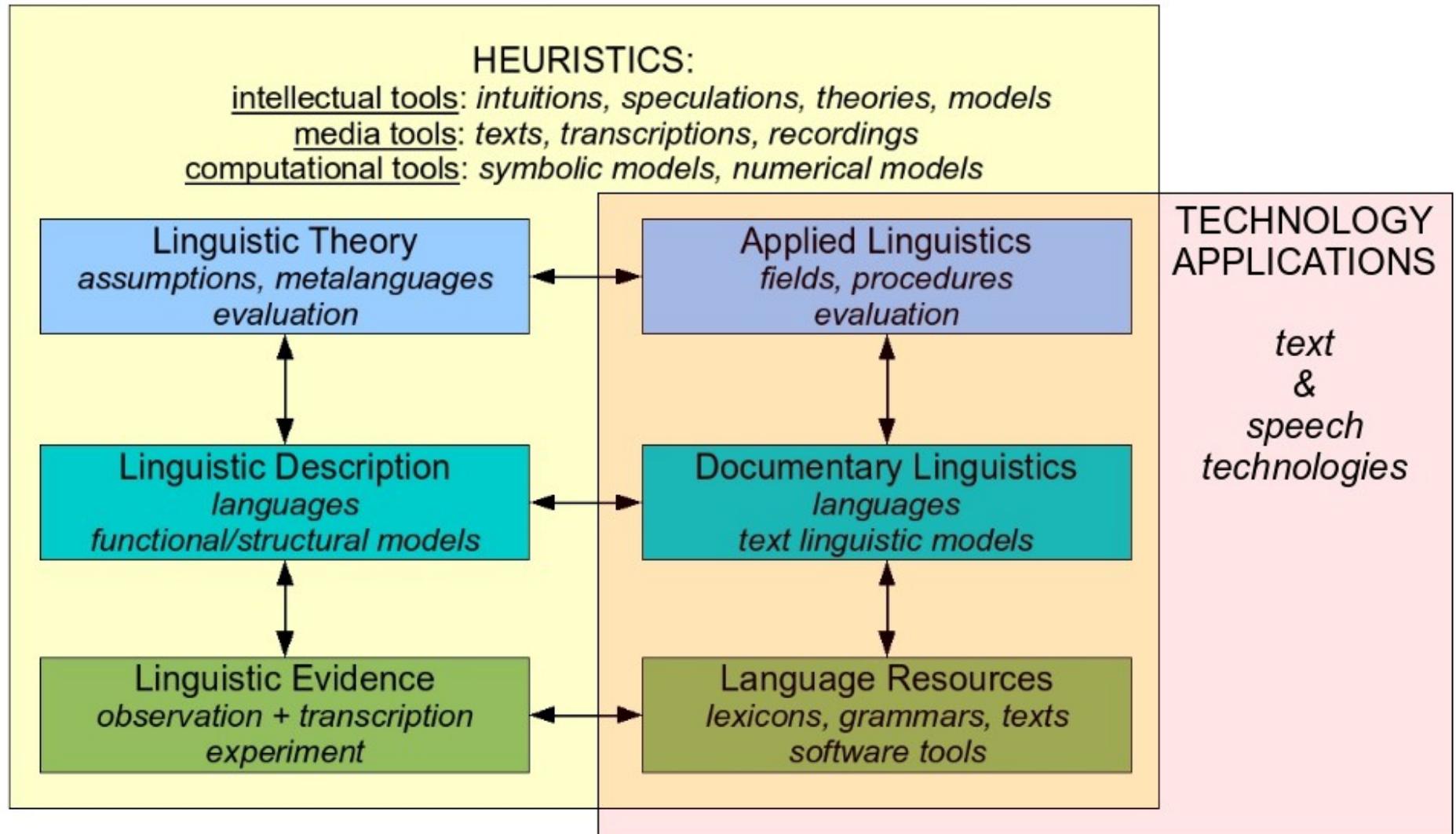
# Language Documentation Workflow (Thieberger)



# Language Documentation: simplified workflow



# Interdisciplinary cooperation



# Language documentation cuisine

<b>Genre:</b> Crêpes	<b>Plan:</b> Difficulté: Facile	<b>Data:</b> 250g de farine
	Préparation: 10 mn	4 oeufs
	Cuisson: 15 mn	1/2l de lait
	Temps Total: 25 mn	1 pincée de sel
	Quantité: 4 pers	2 cuillères à soupe de sucre
		50 g de beurre fondu

## Production:

Mettre la farine dans un **saladier** avec le sel et le sucre. Faire un puits au milieu et y verser les oeufs légèrement battus à la **fourchette**. Commencer à incorporer doucement la farine avec une **cuillère en bois**. Quand le mélange devient épais, ajouter le lait froid petit à petit, on peut utiliser un **fouet** mais toujours doucement pour éviter les grumeaux. Quand tout le lait est mélangé, la pâte doit être assez fluide, si elle vous paraît trop épaisse, rajouter un peu de lait, ensuite le beurre fondu, mélanger bien.

Cuire les crêpes dans une **poêle** chaude (pas besoin de matière grasse, elle est déjà dans la pâte). Verser une petite louche de pâte dans la **poêle**, faire un mouvement de rotation pour répartir la pâte sur toute la surface, poser sur le **feu** et quand le tour de la crêpe se colore en roux clair, il est temps de la retourner. Laisser cuire environ une minute de ce côté et la crêpe est prête. Répéter jusqu'à épuisement de la pâte.

**Evaluation:** Cuisinez, savourez... puis si vous le souhaitez, partagez / déposez (ci-dessous) votre avis sur cette recette.

# Language documentation cuisine

**Genre:** Crêpes

**Plan:** Difficulté: Facile  
Préparation: 10 mn  
Cuisson: 15 mn  
Temps Total: 25 mn  
Quantité: 4 pers

**Data:** 250g de farine  
4 oeufs  
1/2l de lait  
1 pincée de sel  
2 cuillères à soupe de sucre  
50 g de beurre fondu

## Production:

Mettre la farine dans un **saladier** avec le sel et le sucre. Faire un puits au milieu et y verser les oeufs légèrement battus à la **fourchette**. Commencer à incorporer doucement la farine avec une **cuillère en bois**. Quand le mélange devient épais, ajouter le lait froid petit à petit, on peut utiliser un **fouet** mais toujours doucement pour éviter les grumeaux. Quand tout le lait est mélangé, la pâte doit être assez fluide, si elle vous paraît trop épaisse, rajouter un peu de lait, ensuite le beurre fondu, mélanger bien.

Cuire les crêpes dans une **poêle** chaude (pas besoin de matière grasse, elle est déjà dans la pâte). Verser une petite louche de pâte dans la **poêle**, faire un mouvement de rotation pour répartir la pâte sur toute la surface, poser sur le **feu** et quand le tour de la crêpe se colore en roux clair, il est temps de la retourner. Laisser cuire environ une minute de ce côté et la crêpe est prête. Répéter jusqu'à épuisement de la pâte.

**Evaluation:** Cuisinez, savourez... puis si vous le souhaitez (ci-dessous) votre avis sur cette recette.

## Application:

Servir avec du sucre ou du fromage sur une assiette chaude.

# Language documentation cuisine

**Genre:** Crêpes

**Plan:** Difficulté: Facile  
Préparation: 10 mn  
Cuisson: 15 mn  
Temps Total: 25 mn  
Quantité: 4 pers

**Data:** 250g de farine  
4 oeufs  
1/2l de lait  
1 pincée de sel  
2 cuillères à soupe de sucre  
50 g de beurre fondu

## Production:

Mettre la farine dans un **saladier** avec le sel et le sucre. Faire un puits au milieu et y verser les oeufs légèrement battus à la **fourchette**. Commencer à incorporer doucement la farine avec une **cuillère en bois**. Quand le mélange devient épais, ajouter le lait froid petit à petit, on peut utiliser un **fouet** mais toujours doucement pour éviter les grumeaux. Quand tout le lait est mélangé, la pâte doit être assez fluide, si elle vous paraît trop épaisse, rajouter un peu de lait, ensuite le beurre fondu, mélanger bien.

Cuire les crêpes sur une **poêle** chaude (pas besoin de matière grasse, elle est déjà dans la pâte). Verser une louche de pâte dans la **poêle**, faire un mouvement de rotation pour répartir la pâte sur toute la surface, poser sur le **feu** et quand le tour de la crêpe se colore en roux clair, il est temps de la retourner. Laisser cuire environ une minute de ce côté et la crêpe est prête. Répéter jusqu'à épuisement de la pâte.

**Tools**

**Evaluation:** Cuisinez, savourez... puis si vous le souhaitez (ci-dessous) votre avis sur cette recette.

## Application:

Servir avec du sucre ou du fromage sur une assiette chaude.

# What is language documentation?

There are many approaches to language documentation:

- field linguistics
- laboratory linguistics
- armchair linguistics
- 
- strongly empiricist
- qualitative and quantitative methods
- observations rather than generalisations
- digital and multimedia products:
  - transcriptions and texts
  - dictionaries
  - grammars
  - audio: recordings of many different data types
  - video: photos, diagrammes, maps, videos

# **Text linguistics**

The linguistic study of texts is not too different from the linguistic study of sentences or words:

- text syntax:
  - the arrangement of text objects, e.g. of sentences in paragraphs and of paragraphs, pictures, tables etc. in documents
- text semantics:
  - sense:
    - argumentation, description, narration, ...
  - reference:
    - the domain of the text, i.e. languages and language sets
- text pragmatics:
  - intention of writer, design in relation to reader

# Text linguistics

Linguistic models of texts are not too different from models of sentences or words.

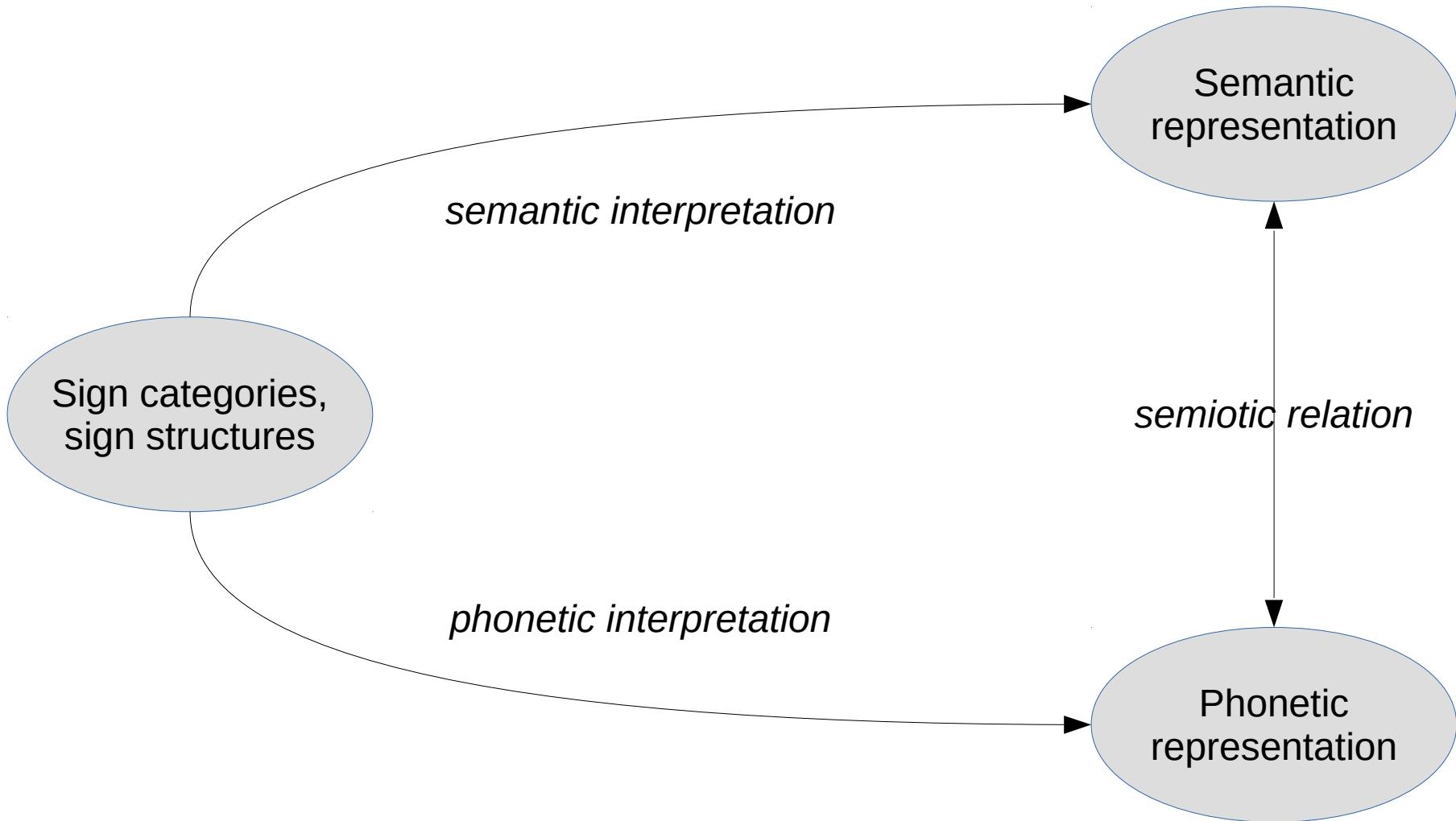
One conventional model:

## syntax of texts

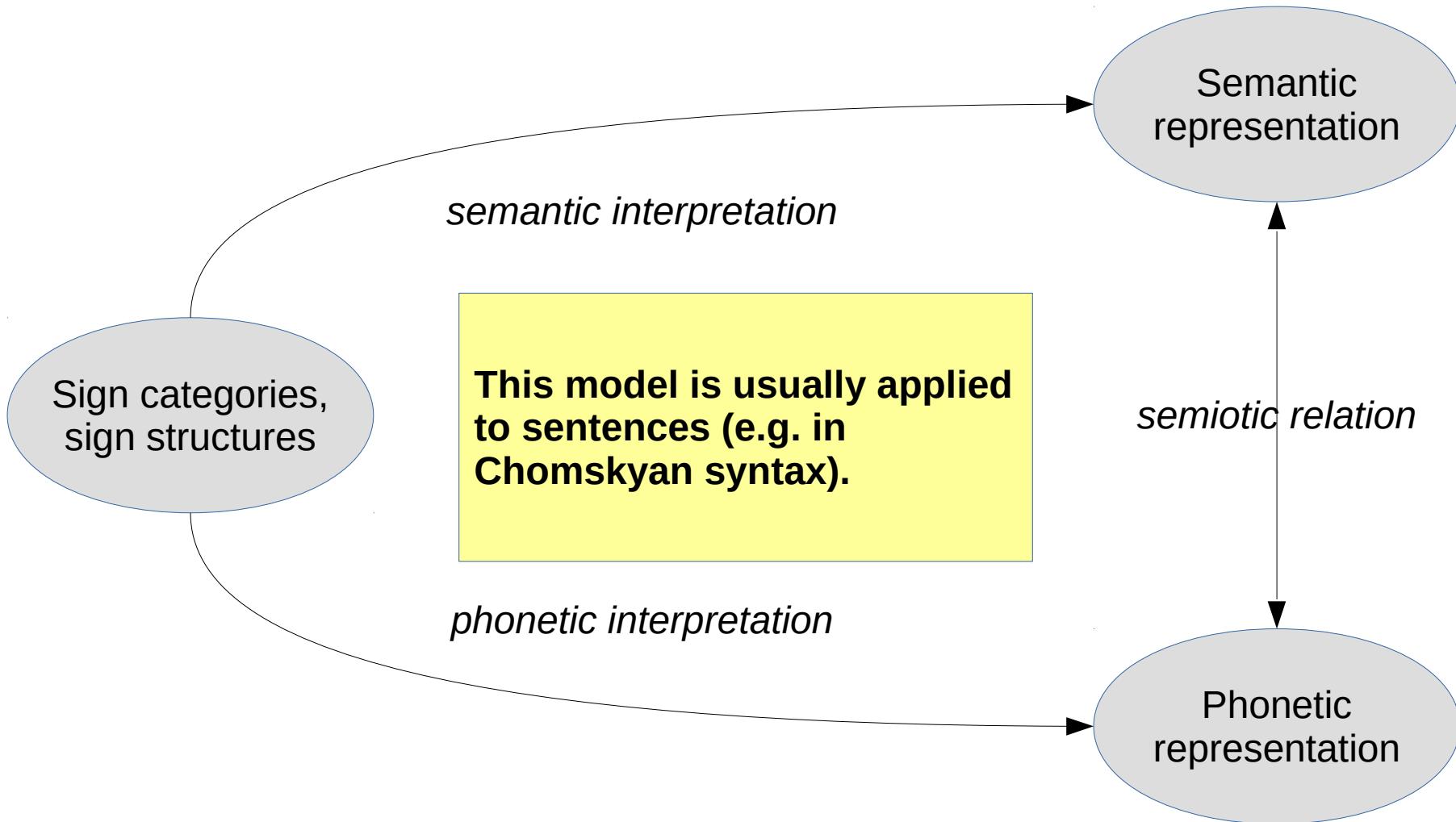
- semantic interpretation of texts
  - sense, reference (previous slide)
- phonetic interpretation of texts
  - prosody, rhythm, voice quality

But what about writing? Multimodal objects?

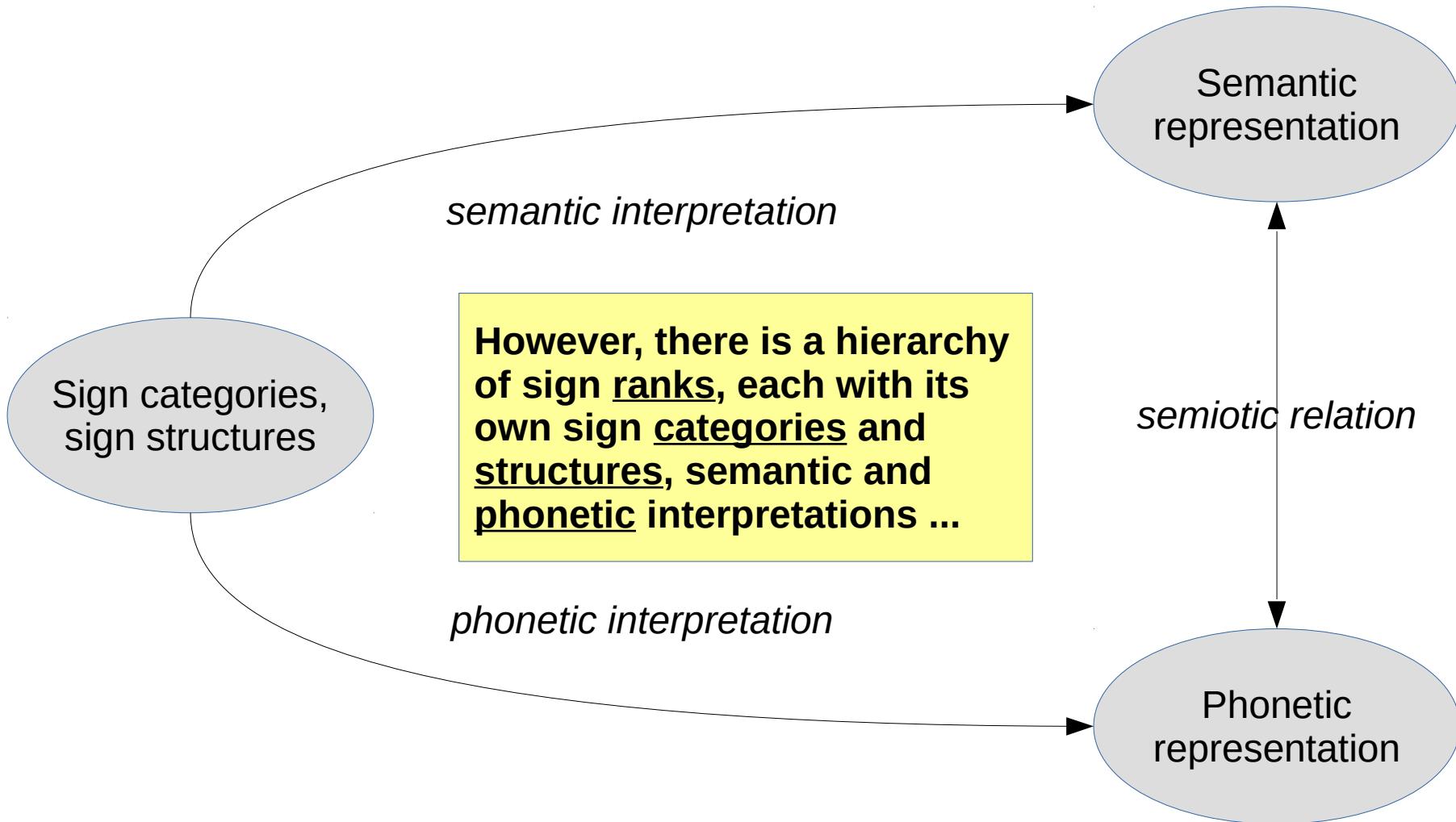
# Text linguistics: a general sign model



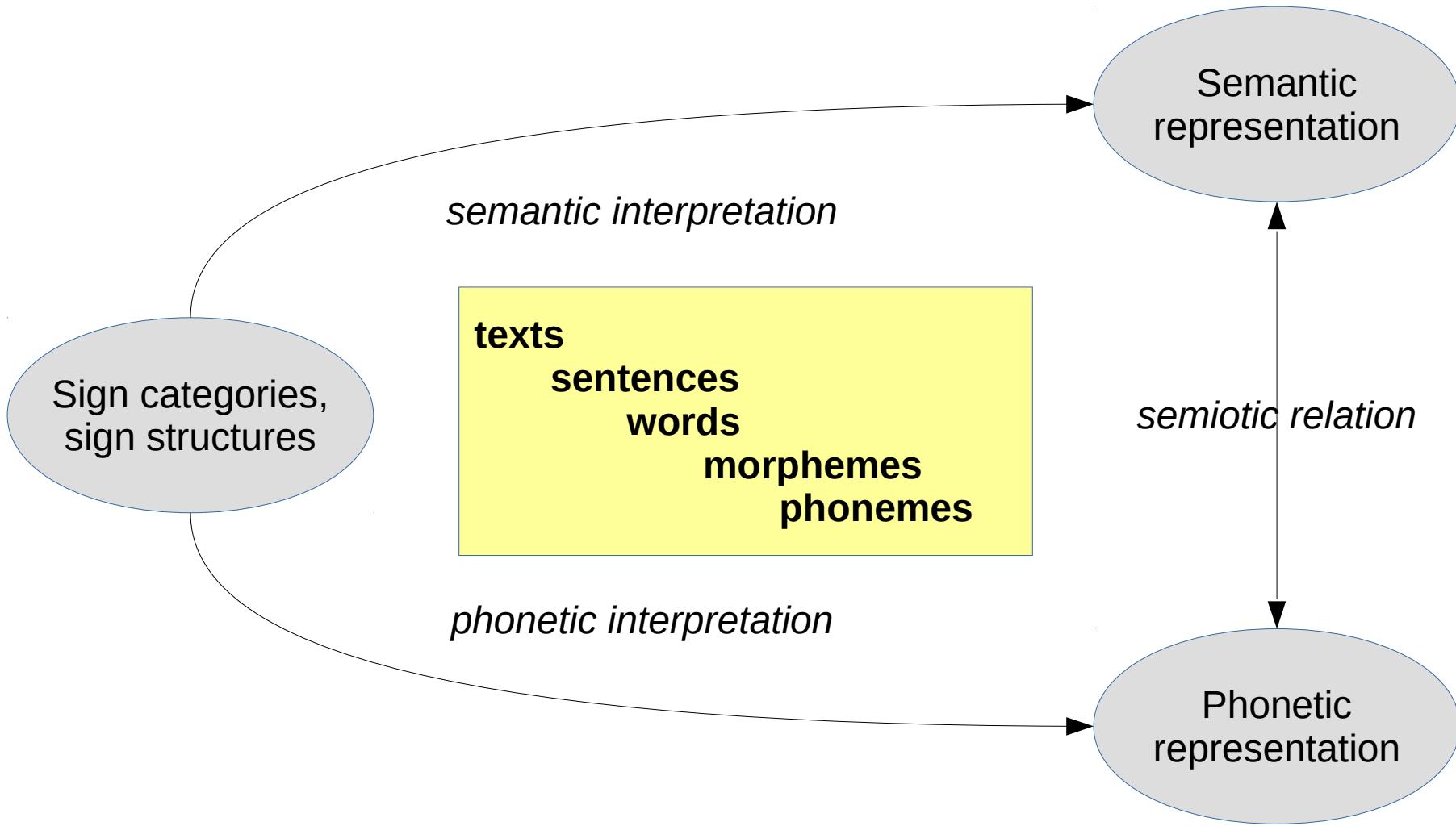
# Text linguistics: a general sign model



# Text linguistics: a general sign model



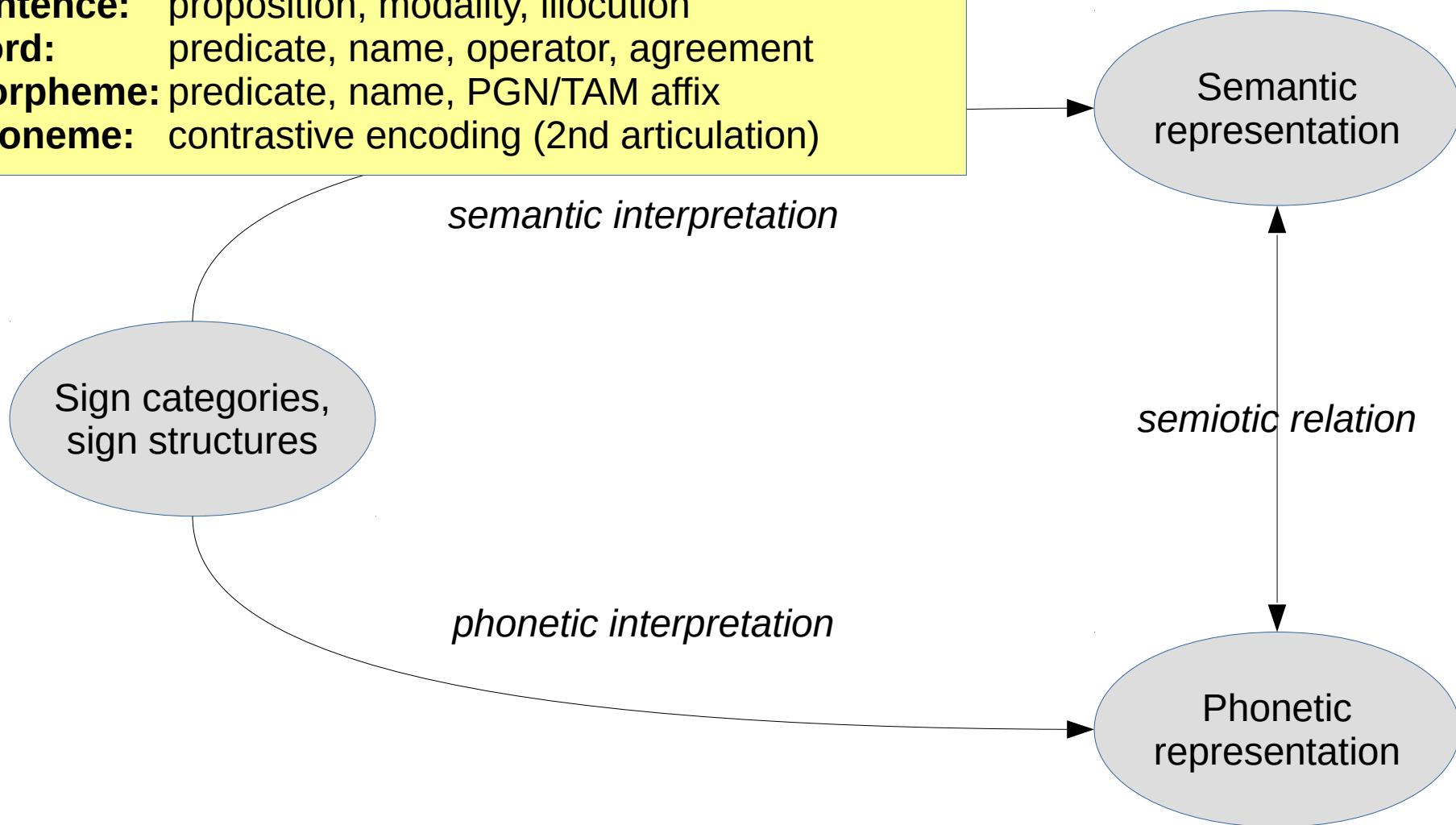
# Text linguistics: a general sign model



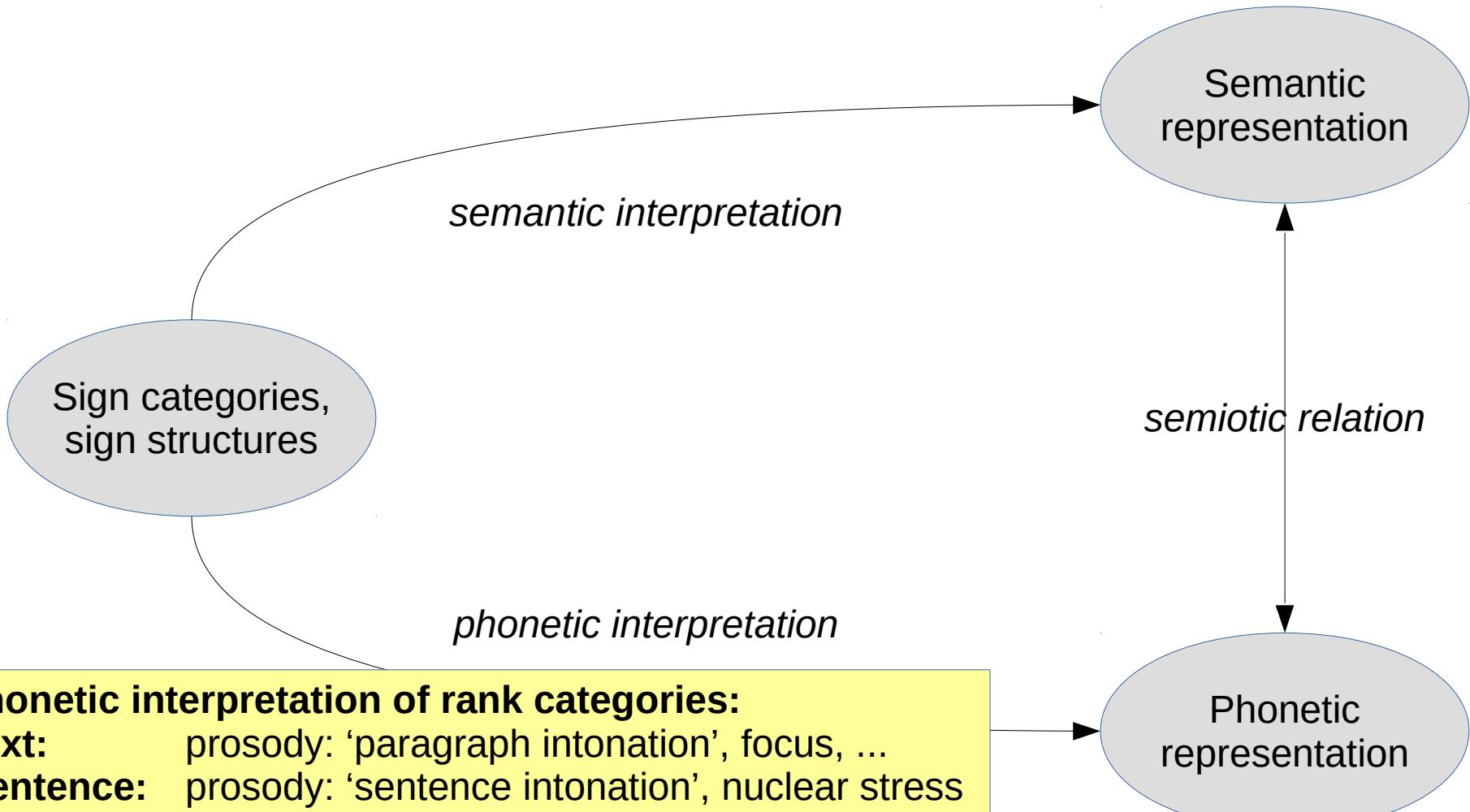
# Text linguistics: a general sign model

## Semantic interpretation of rank categories:

- text:** narrative, argument, illustration, ...
- sentence:** proposition, modality, illocution
- word:** predicate, name, operator, agreement
- morpheme:** predicate, name, PGN/TAM affix
- phoneme:** contrastive encoding (2nd articulation)



# Text linguistics: a general sign model



## Phonetic interpretation of rank categories:

**text:** prosody: 'paragraph intonation', focus, ...

**sentence:** prosody: 'sentence intonation', nuclear stress

**word:** prosody: word stress, e.g. compound stress

**morpheme:** prosody: pitch accent, tone

**phoneme:** distinctive features, tone

# Text linguistics: a general sign model

## Semantic interpretation of rank categories:

- text:** narrative, argument, illustration, ...
- sentence:** proposition, modality, illocution
- word:** predicate, name, operator, agreement
- morpheme:** predicate, name, PGN/TAM affix
- phoneme:** contrastive encoding (2nd articulation)

*semantic interpretation*

Sign categories,  
sign structures

Semantic  
representation

*phonetic interpretation*

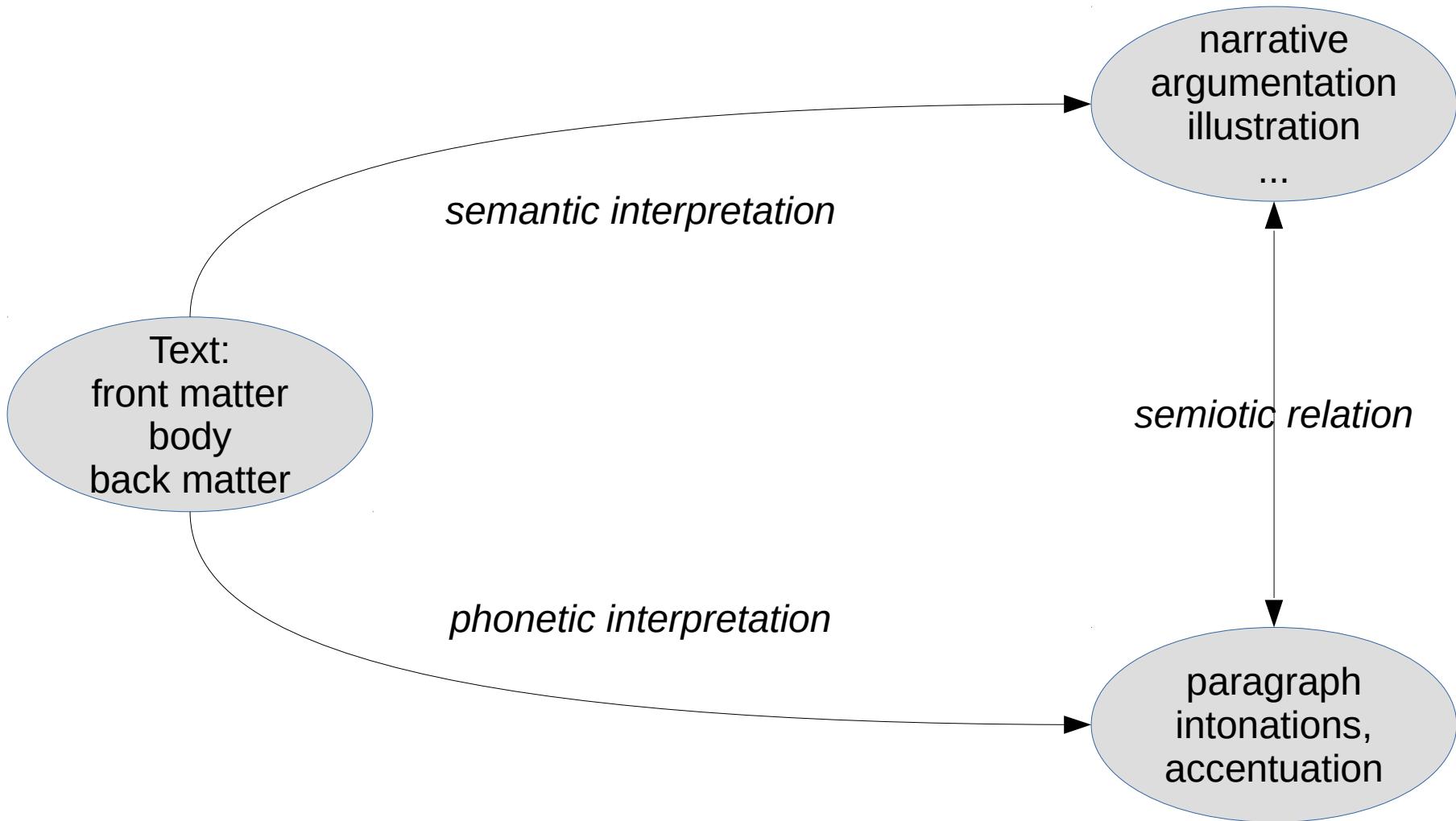
## Phonetic interpretation of rank categories:

- text:** prosody: 'paragraph intonation', focus, ...
- sentence:** prosody: 'sentence intonation', nuclear stress
- word:** prosody: word stress, e.g. compound stress
- morpheme:** prosody: pitch accent, tone
- phoneme:** distinctive features, tone

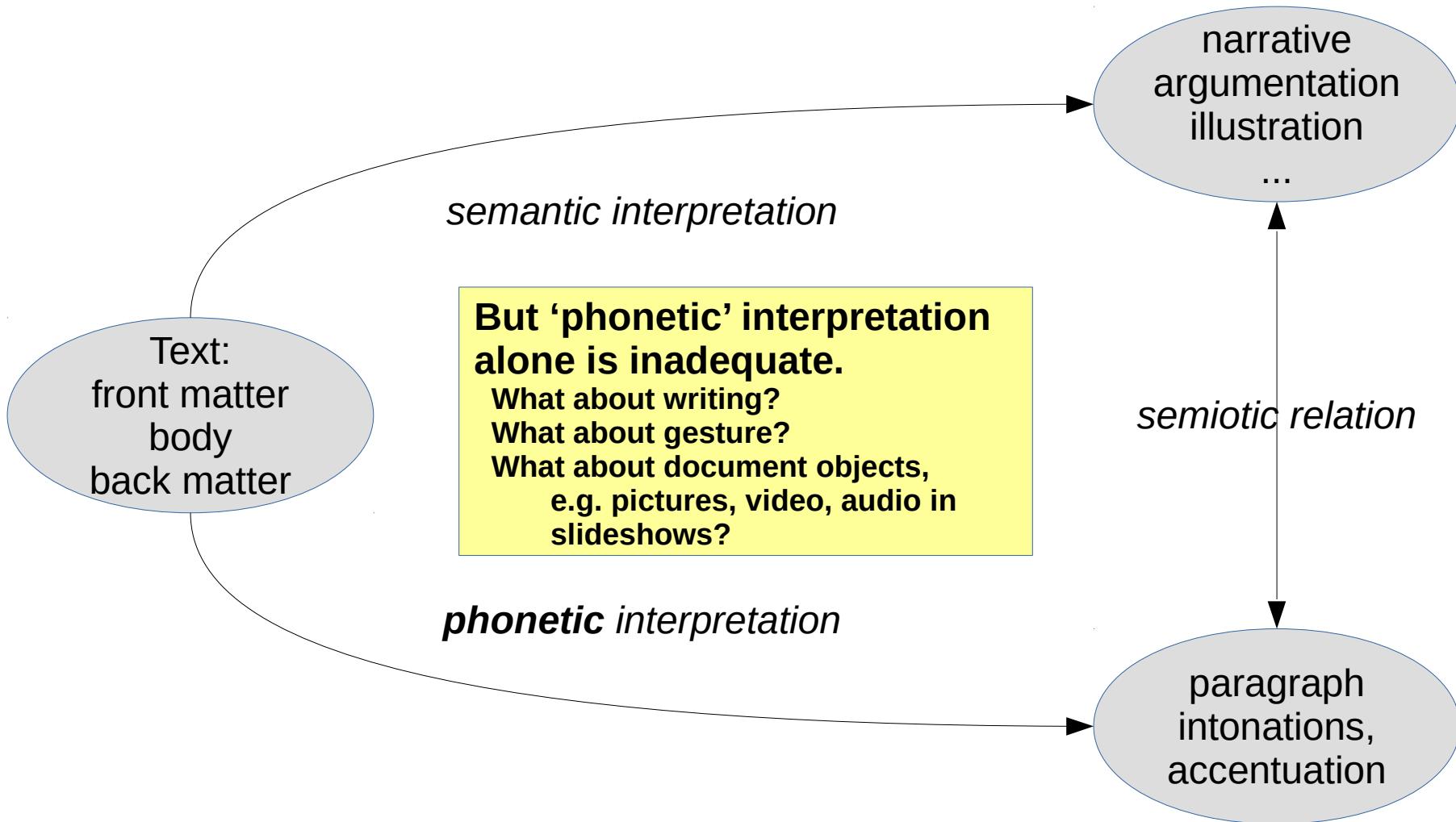
*semiotic relation*

Phonetic  
representation

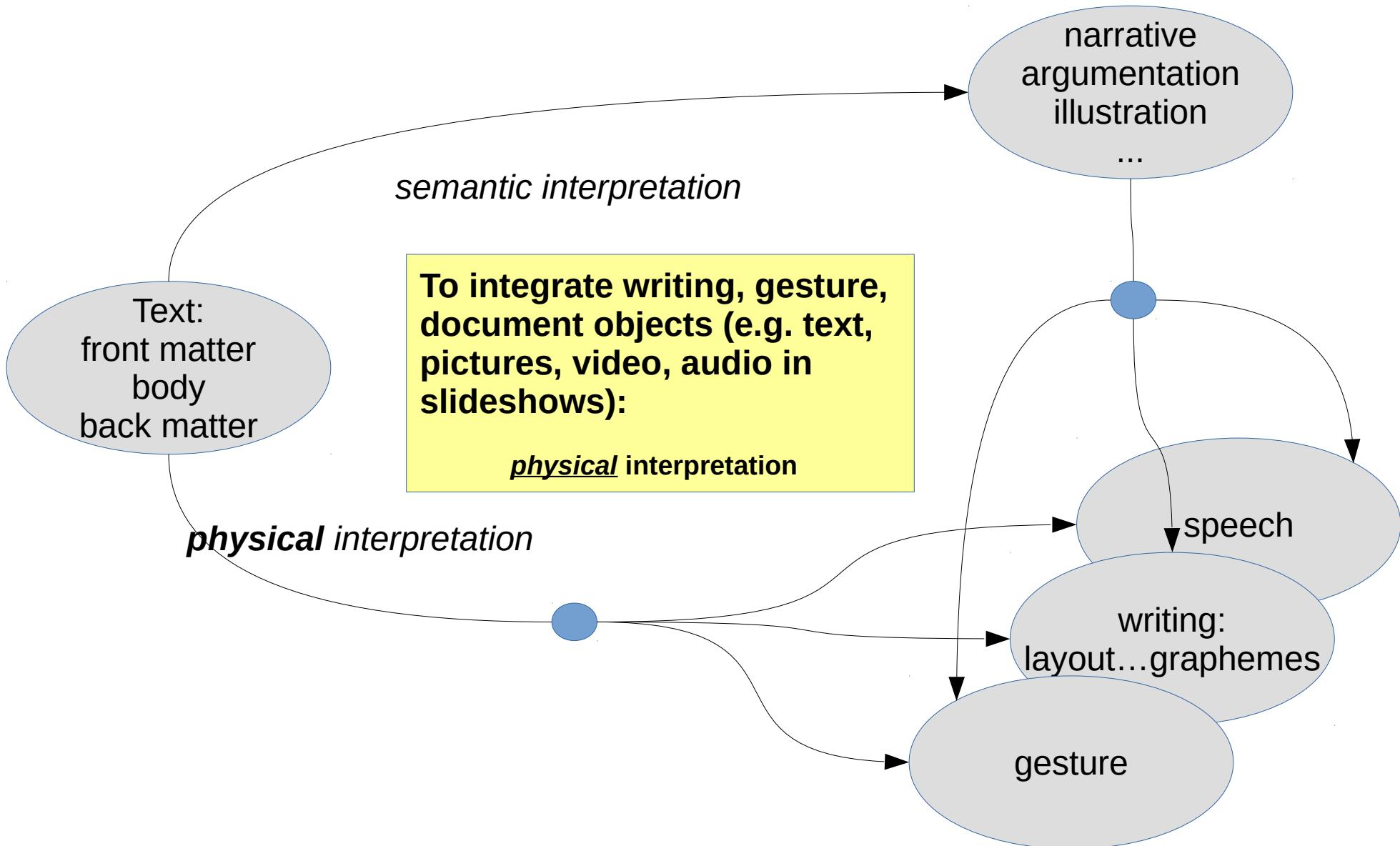
# Text linguistics: a general sign model



# Text linguistics: a general sign model



# Text linguistics: a general sign model



# Text linguistics: a general sign model

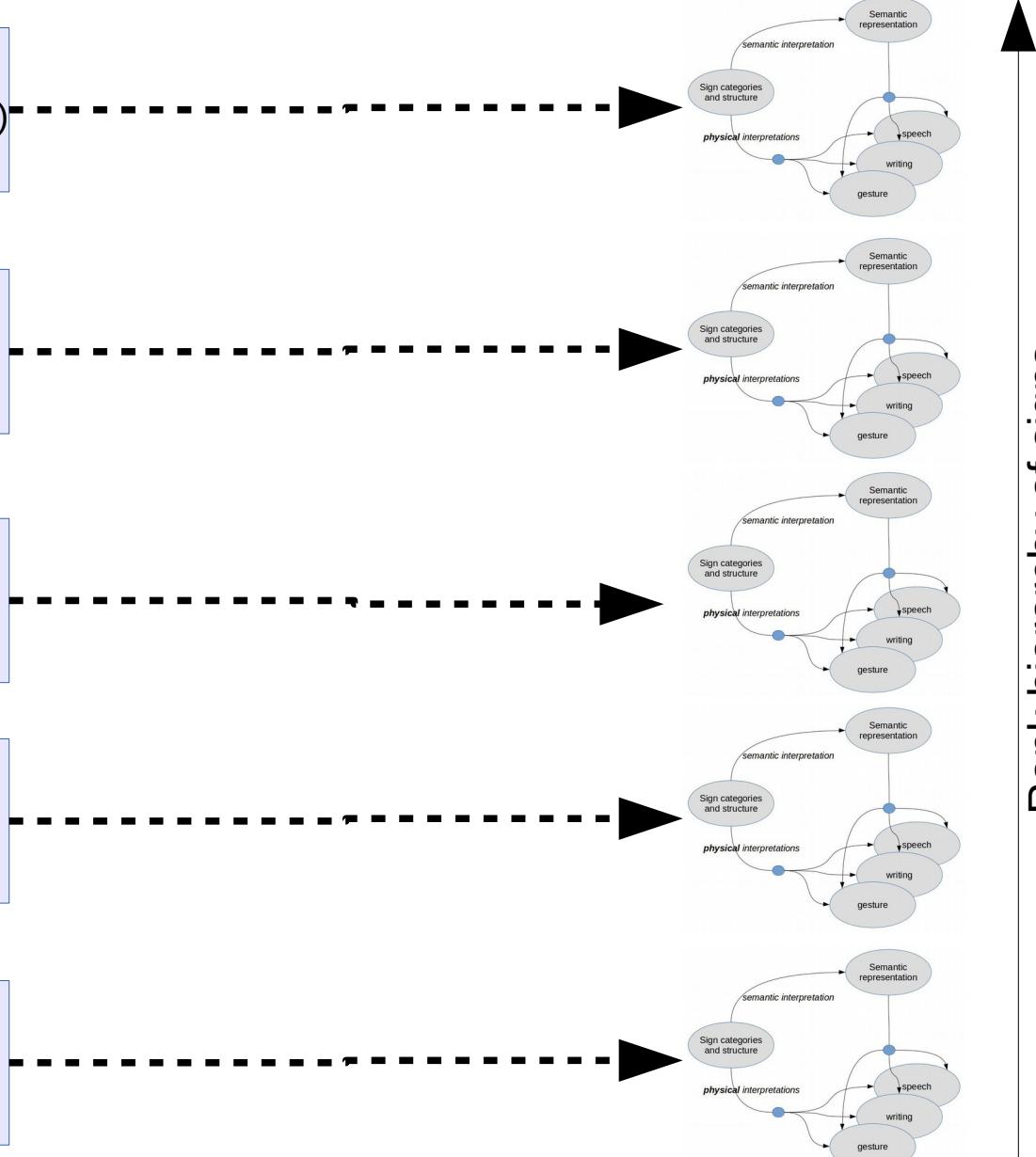
**Text objects**  
(sections, paragraphs, ...) and their properties

**Sentence objects**  
(categories, functions) and their properties

**Word objects**  
and their properties

**Morpheme objects**  
and their properties

**Phoneme objects**  
and their properties



# Analysis of a sample document

**HOW NATURAL IS CHINESE L2 ENGLISH PROSODY?**  
 Joe Yu  
 Dafydd Gibbon  
 Universitat Bielefeld,  
 Bielefeld, Germany  
 gibbon@uni-bielefeld.de

## ABSTRACT

Standard varieties of Chinese and English have many typological similarities for Chinese L2 learners of English at all levels: first, differences in the phonetic system; second, differences in syllable sequence patterns; third, differences in lexical stress-accent language; fourth, timing differences in the prosodic hierarchy, including the timing of intonation, stress, rhythm, and tone. The English-Chinese teaching materials in Chinese L2 and English native speakers are in respect of phonology and phonotactics, and the phonology interface. The SPASS and TGA phonetic analysis tools are used. Results indicate that native speakers are more advanced than L2 proficiency levels and native patients.

**Keywords:** Timing, grammar, L2, Chinese, English.

**1. OBJECTIVES AND BACKGROUND**  
 Standard varieties of Chinese and English have several major typological prosodic differences. These differences are reflected in the phonetic analysis results for Chinese L2 learners of English native varieties of English: (1) different intonation contours; (2) different syllable and syllable sequence patterns; (3) the difference in the timing of intonation and tone; (4) the difference in the timing of stress in the language; (5) timing differences at all levels in the prosodic hierarchy, including differences in prosodic duration, relative stress, etc. [1-4].  
**1.1. METHODS AND DATA**  
 1.1.1. Methods  
 Speech recordings of Chinese L2 speakers and English native speakers were automatically segmented and labeled by the TGA system, and then saved in standard Praat long format [4]. The data were analyzed in two ways. First, we investigated for temporal properties and temporal structures using an online tool which provides heuristics for automatically investigating

temporal properties and distributions in measured data, including global, local and structural measures for duration, rate and tempo, TGA's online tool of [3]. English proficiency levels of the L2 speakers were tested, and temporal properties and distributions of the speech samples were compared with their proficiency levels.

## 1.2. Data

The data used here is used because it is a standard feature of EFL learners. Data are from the English-CASE Chinese English learner corpus [14], which contains both speech samples of native English speakers and speech samples of non-native Chinese speakers. The corpus includes the Aesop's fable *The North Wind and the Sun*.

## 1.3. RESULTS AND DISCUSSION

### 1.3.1. Proficiency evaluation

The pre-evaluation was done by 4 Chinese English teachers. The final evaluation was done by 30 English native speakers graded on a scale of 5 by gender impression: excellent, good, average, poor, bad.

### 1.3.2. Speakers

Speakers are further evaluated on a 5-point scale: native speaker, advanced, medium, poor, native speaker, native speaker, etc. sequences. These impressions are based on the following two properties:

#### (1) temporal pattern distribution as evidence for language learning and grammaticalization.

#### (2) intonation, stress, rhythm etc. There is no necessary correlation between speech rate and intonation.

The primary objective is to obtain new findings on the differences in prosody between Chinese L2 English and English native speakers.

The practical objective is to provide a basis for creating general guidelines on timing, intonation, stress, rhythm etc. for teaching English, pronunciation and spelling.

To test validity, the final evaluation results are compared between and within Chinese English and native English speakers, and the results are combined.

The native English speakers had higher scores for prosodic criteria than English teachers, indicating either slurred and L2 features, or more natural and native-like intonation.

Table 1 shows the proficiency of the female learners to be lower than that of the male.

$F(1,18) = 4.18, p < 0.05$ . Language proficiency does not correlate with years of learning,  $r^2 = 0.214, p > 0.05$ .

Table 1: Chinese learner proficiency in English.

	English proficiency	advanced	medium	poor
Gender	male	3	4	5
	female	1	0	2
Learning (years)	10	1	5	5
	20	1	5	5

Native speakers are more advanced than non-native speakers in terms of English proficiency.

## 1.3.2. Temporal dispersion: Wagner Quadrant

A more informative technique than the older global metrics is the Wagner Quadrant method [11]. This technique is based on the analysis of the distribution of the duration of the n-gran. Figure 3. The scatter plots are 2-scores of duration of the n-gran. The four quadrants of the quadrant of the plots around zero is  $x=2 \times \text{shorter-changer} - 2 \times \text{longer-changer}$ ;  $y = 2 \times \text{longer-changer} - 2 \times \text{shorter-changer}$ . The most advanced non-native Chinese reader (Figure 1) will show a Wagner Quadrant with a large area of the quadrant below that of the advanced Chinese reader (Figure 2) to the right of the quadrant below.

Vocal inspection of the sample figures shows that that tendency is indeed present: the low proficiency speakers have a much larger area of values through the four quadrants.

The distribution of the n-gran, on the other hand, tends to cluster values in the shorter-longer and longer-shorter quadrant, reflecting the tendency of English native speakers to have a distribution which showed a strong syllable alternating with unvoiced consonants and vowels and short and short, respectively. There are many *shorter-changer* and *longer-changer* patterns in the early patterns of a long syllable followed by more than one short syllables. The distribution is apparently more uniform in Chinese. This is also evident in the case of the advanced Chinese speakers, who have a more uniform pattern as predicted by our initial hypothesis.

Table 2: Summary of mean variability and mean syllable rate.

	Chinese L1	Chinese L2	Chinese L2	English
poor	4.5	8.5	9.5	13.1
medium	4.9	7.5	7.5	12.1
advanced	5.1	2.8	-	12.2
native	5.2	2.4	-	9.8

Quantitative analysis of each of the four quadrants for all speakers is planned for future work.

## 1.3.3. Quantitative analysis of each of the four quadrants for all speakers is planned for future work.

### 1.3.4. Temporal n-gran

Traditional methods for quantifying the timing of speech provide single global indices of mean duration, mean syllable rate, etc. However, there are differences in the structure of the speech signal, and these differences have been developed solely to examine the sequential structure of binary alternation patterns in syllable sequences. The traditional methods of quantifying rhythm, relative accentedness, covered by the global measures of duration and syllable rate, do not measure the duration difference n-gran as categorical, at this stage. Intensity and other rhythm patterns are not dealt with.

## 1.3.5. Inexperienced reader of English overrides the syllable near-synchrony of the native language

Figure 2: Chinese L2 English advanced, female.

Figure 2: Chinese L2 English advanced, female.

Figure 3: Wagner Quadrant for female native USA.

Figure 3: Wagner Quadrant for female native USA.

Figure 4: Quantitative analysis of each of the four quadrants for all speakers is planned for future work.

Figure 5: Time Tree.

Figure 6: Time Tree.

Figure 7: Time Tree.

Figure 8: Time Tree.

Figure 9: Time Tree.

Figure 10: Time Tree.

Figure 11: Time Tree.

Figure 12: Time Tree.

Figure 13: Time Tree.

Figure 14: Time Tree.

Figure 15: Time Tree.

Figure 16: Time Tree.

Figure 17: Time Tree.

Figure 18: Time Tree.

Figure 19: Time Tree.

Figure 20: Time Tree.

Figure 21: Time Tree.

Figure 22: Time Tree.

Figure 23: Time Tree.

Figure 24: Time Tree.

Figure 25: Time Tree.

Figure 26: Time Tree.

Figure 27: Time Tree.

Figure 28: Time Tree.

Figure 29: Time Tree.

Figure 30: Time Tree.

Figure 31: Time Tree.

Figure 32: Time Tree.

Figure 33: Time Tree.

Figure 34: Time Tree.

Figure 35: Time Tree.

Figure 36: Time Tree.

Figure 37: Time Tree.

Figure 38: Time Tree.

Figure 39: Time Tree.

Figure 40: Time Tree.

Figure 41: Time Tree.

Figure 42: Time Tree.

Figure 43: Time Tree.

Figure 44: Time Tree.

Figure 45: Time Tree.

Figure 46: Time Tree.

Figure 47: Time Tree.

Figure 48: Time Tree.

Figure 49: Time Tree.

Figure 50: Time Tree.

Figure 51: Time Tree.

Figure 52: Time Tree.

Figure 53: Time Tree.

Figure 54: Time Tree.

Figure 55: Time Tree.

Figure 56: Time Tree.

Figure 57: Time Tree.

Figure 58: Time Tree.

Figure 59: Time Tree.

Figure 60: Time Tree.

Figure 61: Time Tree.

Figure 62: Time Tree.

Figure 63: Time Tree.

Figure 64: Time Tree.

Figure 65: Time Tree.

Figure 66: Time Tree.

Figure 67: Time Tree.

Figure 68: Time Tree.

Figure 69: Time Tree.

Figure 70: Time Tree.

Figure 71: Time Tree.

Figure 72: Time Tree.

Figure 73: Time Tree.

Figure 74: Time Tree.

Figure 75: Time Tree.

Figure 76: Time Tree.

Figure 77: Time Tree.

Figure 78: Time Tree.

Figure 79: Time Tree.

Figure 80: Time Tree.

Figure 81: Time Tree.

Figure 82: Time Tree.

Figure 83: Time Tree.

Figure 84: Time Tree.

Figure 85: Time Tree.

Figure 86: Time Tree.

Figure 87: Time Tree.

Figure 88: Time Tree.

Figure 89: Time Tree.

Figure 90: Time Tree.

Figure 91: Time Tree.

Figure 92: Time Tree.

Figure 93: Time Tree.

Figure 94: Time Tree.

Figure 95: Time Tree.

Figure 96: Time Tree.

Figure 97: Time Tree.

Figure 98: Time Tree.

Figure 99: Time Tree.

Figure 100: Time Tree.

Figure 101: Time Tree.

Figure 102: Time Tree.

Figure 103: Time Tree.

Figure 104: Time Tree.

Figure 105: Time Tree.

Figure 106: Time Tree.

Figure 107: Time Tree.

Figure 108: Time Tree.

Figure 109: Time Tree.

Figure 110: Time Tree.

Figure 111: Time Tree.

Figure 112: Time Tree.

Figure 113: Time Tree.

Figure 114: Time Tree.

Figure 115: Time Tree.

Figure 116: Time Tree.

Figure 117: Time Tree.

Figure 118: Time Tree.

Figure 119: Time Tree.

Figure 120: Time Tree.

Figure 121: Time Tree.

Figure 122: Time Tree.

Figure 123: Time Tree.

Figure 124: Time Tree.

Figure 125: Time Tree.

Figure 126: Time Tree.

Figure 127: Time Tree.

Figure 128: Time Tree.

Figure 129: Time Tree.

Figure 130: Time Tree.

Figure 131: Time Tree.

Figure 132: Time Tree.

Figure 133: Time Tree.

Figure 134: Time Tree.

Figure 135: Time Tree.

Figure 136: Time Tree.

Figure 137: Time Tree.

Figure 138: Time Tree.

Figure 139: Time Tree.

Figure 140: Time Tree.

Figure 141: Time Tree.

Figure 142: Time Tree.

Figure 143: Time Tree.

Figure 144: Time Tree.

Figure 145: Time Tree.

Figure 146: Time Tree.

Figure 147: Time Tree.

Figure 148: Time Tree.

Figure 149: Time Tree.

Figure 150: Time Tree.

Figure 151: Time Tree.

Figure 152: Time Tree.

Figure 153: Time Tree.

Figure 154: Time Tree.

Figure 155: Time Tree.

Figure 156: Time Tree.

Figure 157: Time Tree.

Figure 158: Time Tree.

Figure 159: Time Tree.

Figure 160: Time Tree.

Figure 161: Time Tree.

Figure 162: Time Tree.

Figure 163: Time Tree.

Figure 164: Time Tree.

Figure 165: Time Tree.

Figure 166: Time Tree.

Figure 167: Time Tree.

Figure 168: Time Tree.

Figure 169: Time Tree.

Figure 170: Time Tree.

Figure 171: Time Tree.

Figure 172: Time Tree.

Figure 173: Time Tree.

Figure 174: Time Tree.

Figure 175: Time Tree.

Figure 176: Time Tree.

Figure 177: Time Tree.

Figure 178: Time Tree.

Figure 179: Time Tree.

Figure 180: Time Tree.

# Goal of this course

The practical goal of this course is to create

- a professional quality report
- with a structure of this kind
- to document your chosen language
- using whatever document object and properties are necessary

The document objects will be:

- title, author, affiliation
- table of contents
- abstract
- sections and subsections with headings
- paragraphs
- tables (e.g. dictionary)
- figures (e.g. map)
- references

# Document objects and properties

The document objects will be:

- title, author, affiliation
- table of contents
- abstract
- sections and subsections with headings
- paragraphs
- tables (e.g. dictionary)
- figures (e.g. map)
- references

The properties of the objects will be:

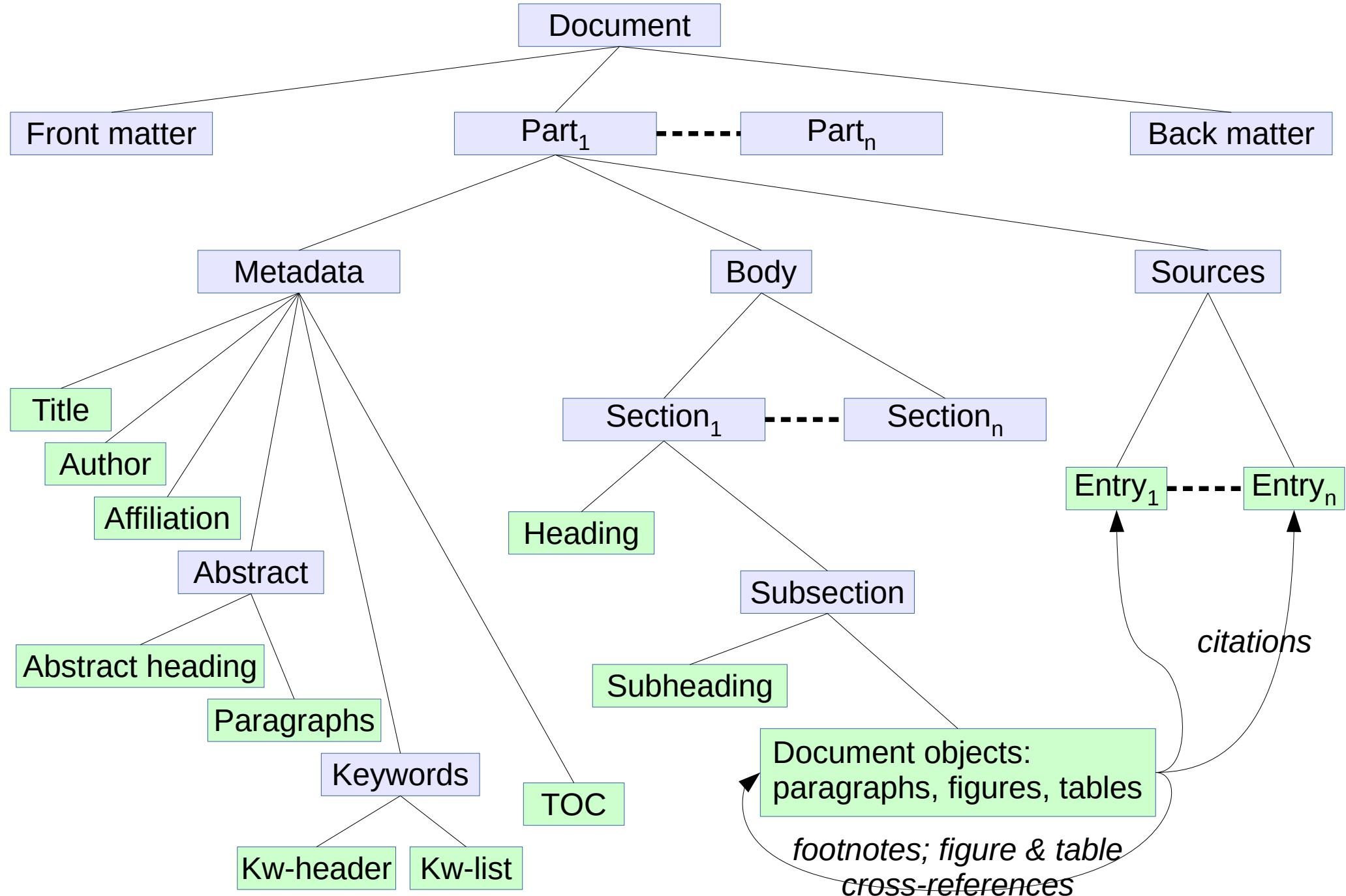
- separation from preceding and following objects
- alignment (left, right, centre, justified)
- for tables and figures: caption
- for paragraphs: indentation, format style

# Typical structure of a dictionary: the ‘5M’ model

## Megastructure (the entire document):

- Metastructure (about the document)
  - front matter
  - sketch grammar
  - abbreviations
  - back matter
- Macrostructure (e.g. semasiological, onomasiological)
  - Entries
- Microstructure (of entries)
  - types of lexical information (data categories)
- Mesostructure (wordnet structure)
  - cross references between entries (e.g. synonyms, antonyms)

# **Structure of a typical report document or article**



# Practical document creation: styles are text grammar

Topic: Documentation of a language

Method: Word processing

- OpenOffice
- LibreOffice
- MS Word
- LaTeX
- ...
- Next meetings:
  - other document types
  - databases

# Practical document creation: styles are text grammar

Topic: a language

Method: Word processing

Documentation tool:LibreOffice

Specification:

Document syntax:

- as previously illustrated, with tables & pictures

Semantics:

- Information about a language from the courses, the internet, your knowledge, ...

Physical interpretation:

- screen, paper, internet

**Time for practical work ...**