# Of Nodes and Cells.

## Two Perspectives on (and from) *Word Formation Latin*

**Marco Budassi**, **Eleonora Litta** and **Marco Passarotti**

ICLL 2019 | Universidad de Las Palmas de Gran Canaria | 19 June 2019

WFL: Word formation-based lexicon for Classical Latin

WFL: Word formation-based lexicon for Classical Latin

► Lexical base => **Lemlat**: morphological analyser and lemmatiser for Classical and Late Latin (OLD, Georges, Gradenwitz)

WFL: Word formation-based lexicon for Classical Latin

- ► Lexical base => **Lemlat**: morphological analyser and lemmatiser for Classical and Late Latin (OLD, Georges, Gradenwitz)
- ► WFRs are modelled as **directed one-to-many input-output** relations between lemmas (based on I&A models of grammatical description)

WFL: Word formation-based lexicon for Classical Latin

- ► Lexical base => **Lemlat**: morphological analyser and lemmatiser for Classical and Late Latin (OLD, Georges, Gradenwitz)
- ► WFRs are modelled as **directed one-to-many input-output** relations between lemmas (based on I&A models of grammatical description)
- ► **Morphotactic** approach: each WF process is treated individually as the application of one single rule in a certain order.

▶ Relationships between lemmas of the same "word formation family" are represented as the edges in a **directed graph** with a **hierarchical tree-like** structure

- ▶ Relationships between lemmas of the same "word formation family" are represented as the edges in a **directed graph** with a **hierarchical tree-like** structure
- ▶ A **node** is a lemma, and an **edge** is the WFR used to derive the output lemma from the input one, together with any affix

**Compounding** is also shown as an intersection between word formation families.

But: **directed graphs** are not completely satisfactory in representing the full range of relationships included within a word formation family.

But: **directed graphs** are not completely satisfactory in representing the full range of relationships included within a word formation family.
Main problems:

But: **directed graphs** are not completely satisfactory in representing the full range of relationships included within a word formation family.
Main problems:

▶ Directionality

But: **directed graphs** are not completely satisfactory in representing the full range of relationships included within a word formation family.
Main problems:

► Directionality

► Non-linear derivations

The appeal of a **derivational paradigm** is the possibility of stacking recurring word formations into a model, carrying descriptive and predictive powers.

The appeal of a **derivational paradigm** is the possibility of stacking recurring word formations into a model, carrying descriptive and predictive powers.

- ▶ Relationships are not limited to base-derivative pairs (Jackendoff, 1975) => Non linear derivations

The appeal of a **derivational paradigm** is the possibility of stacking recurring word formations into a model, carrying descriptive and predictive powers.

- ► Relationships are not limited to base-derivative pairs (Jackendoff, 1975) => Non linear derivations
- ► No need for directionality (Jackendoff, 1975) => Diachronic discrepancies

The appeal of a **derivational paradigm** is the possibility of stacking recurring word formations into a model, carrying descriptive and predictive powers.

- ► Relationships are not limited to base-derivative pairs (Jackendoff, 1975) => Non linear derivations
- ► No need for directionality (Jackendoff, 1975) => Diachronic discrepancies
- ► Potentiality (Štekauer: 2014) and Availability of Slots more important than the form filling them (Bauer, 1997)

The appeal of a **derivational paradigm** is the possibility of stacking recurring word formations into a model, carrying descriptive and predictive powers.

- ► Relationships are not limited to base-derivative pairs (Jackendoff, 1975) => Non linear derivations
- ► No need for directionality (Jackendoff, 1975) => Diachronic discrepancies
- ► Potentiality (Štekauer: 2014) and Availability of Slots more important than the form filling them (Bauer, 1997)
- ► Regularity and predictability (Bauer, 1997)

The appeal of a **derivational paradigm** is the possibility of stacking recurring word formations into a model, carrying descriptive and predictive powers.

- ► Relationships are not limited to base-derivative pairs (Jackendoff, 1975) => Non linear derivations
- ► No need for directionality (Jackendoff, 1975) => Diachronic discrepancies
- ► Potentiality (Štekauer: 2014) and Availability of Slots more important than the form filling them (Bauer, 1997)
- ► Regularity and predictability (Bauer, 1997)

BUT... not many examples of fuller paradigms (in any language!).

► No directionality: necessary to accommodate those lemmas for which the derivational process is not of the simplex (or simpler) > complex type

- ► No directionality: necessary to accommodate those lemmas for which the derivational process is not of the simplex (or simpler) > complex type
- ► The CELL has a central role in the paradigm (predictability and regularity)

► No directionality: necessary to accommodate those lemmas for which the derivational process is not of the simplex (or simpler) > complex type

► The CELL has a central role in the paradigm (predictability and regularity)

► Each cell must be described in both its morphological characteristics and its semantic features, due to the underlying role of semantics in accounting for derivational processes

The LiLa project: Knowledge Base of Linguistic Resources and Natural Language Processing (NLP) tools for Latin.

The LiLa project: Knowledge Base of Linguistic Resources and Natural Language Processing (NLP) tools for Latin. LiLa collects and connects:

- ▶ Linguistic Resources: corpora, lexica, ontologies, dictionaries, thesauri

The LiLa project: Knowledge Base of Linguistic Resources and Natural Language Processing (NLP) tools for Latin. LiLa collects and connects:

- ▶ Linguistic Resources: corpora, lexica, ontologies, dictionaries, thesauri
- ▶ NLP tools: tokenisers, lemmatisers, PoS-taggers, morphological analysers and dependency parsers

The LiLa project: Knowledge Base of Linguistic Resources and Natural Language Processing (NLP) tools for Latin. LiLa collects and connects:

- ▶ Linguistic Resources: corpora, lexica, ontologies, dictionaries, thesauri
- ▶ NLP tools: tokenisers, lemmatisers, PoS-taggers, morphological analysers and dependency parsers

The Knowledge Base consists of different kinds of objects connected via an explicitly-declared vocabulary for knowledge description.

LiLa is based on an ontology made of:

LiLa is based on an ontology made of:

► **Individuals**: instances of objects (one specific token, lemma etc.)

LiLa is based on an ontology made of:

- ► **Individuals**: instances of objects (one specific token, lemma etc.)
- ► **Classes**: types of objects/concepts (token, lemma, PoS etc.)

LiLa is based on an ontology made of:

- ▶ **Individuals**: instances of objects (one specific token, lemma etc.)
- ▶ **Classes**: types of objects/concepts (token, lemma, PoS etc.)
- ▶ **Attributes**: properties that objects can/must have (morphological features for lemmas/tokens). Each attribute can be a class or an individual

LiLa is based on an ontology made of:

- ▶ **Individuals**: instances of objects (one specific token, lemma etc.)
- ▶ **Classes**: types of objects/concepts (token, lemma, PoS etc.)
- ▶ **Attributes**: properties that objects can/must have (morphological features for lemmas/tokens). Each attribute can be a class or an individual
- ▶ **Relations**: ways in which classes and individuals can be related to one another: RDF triples. Labels from a restricted dictionary of knowledge description: *has_lemma*, *has_PoS*

Different approach to Word Formation:

Different approach to Word Formation:

- ▶ Structure: **declarative** rather than procedural

Different approach to Word Formation:

- ► Structure: **declarative** rather than procedural
- ► No directionality

Different approach to Word Formation:

- ► Structure: **declarative** rather than procedural
- ► No directionality
- ► No morphotaxis.

Most unassuming, declarative way of describing word formation processes.

Most unassuming, declarative way of describing word formation processes.

► Construction: [co(n)] [stell](a)(t)io$_N$ (more specific)

Most unassuming, declarative way of describing word formation processes.

▶ Construction: $[co(n) [stell](a)(t)io]_N$ (more specific)
▶ Schema $[co(n)[x](t)io]_N$ (more generalised)

# Construction Morphology (CxM)

Most unassuming, declarative way of describing word formation processes.

- ▶ Construction: $[co(n) [stell](a)(t)io]_N$ (more specific)
- ▶ Schema $[co(n)[x](t)io]_N$ (more generalised)
- ▶ Constructions and schemas are word-based and declarative

Most unassuming, declarative way of describing word formation processes.

- ► Construction: [co(n) [stell](a)(t)io]$_N$ (more specific)
- ► Schema [co(n)[x](t)io]$_N$ (more generalised)
- ► Constructions and schemas are word-based and declarative
- ► Perfect for LiLa => words are described in their formative elements, which can be organised into connected classes of objects into an ontology.

Three classes of objects:

Three classes of objects:

1. Lemmas

Three classes of objects:

1. Lemmas
2. Affixes (prefixes and suffixes)

Three classes of objects:

1. Lemmas
2. Affixes (prefixes and suffixes)
3. Bases (connectors between lemmas of the same WF family)

Three classes of objects:

1. Lemmas
2. Affixes (prefixes and suffixes)
3. Bases (connectors between lemmas of the same WF family)

Connected by three possibile relationships:

Three classes of objects:

1. Lemmas
2. Affixes (prefixes and suffixes)
3. Bases (connectors between lemmas of the same WF family)

Connected by three possibile relationships:

1. hasPrefix
2. hasSuffix
3. hasBase

LiLa triplestore available at:
https://lila-erc.eu/data/

More sophisticated ways of querying the data:

LiLa triplestore available at:
https://lila-erc.eu/data/

More sophisticated ways of querying the data:

► Start from a lemma (e.g. *formalis*), find its lexical base, retrieve all lemmas and affixes linked to that lexical base (word formation family)

LiLa triplestore available at:
https://lila-erc.eu/data/

More sophisticated ways of querying the data:

► Start from a lemma (e.g. *formalis*), find its lexical base, retrieve all lemmas and affixes linked to that lexical base (word formation family)

► Find lemmas with a specific prefix (ad-) or suffix (-bil)

LiLa triplestore available at:
https://lila-erc.eu/data/

More sophisticated ways of querying the data:

- ► Start from a lemma (e.g. *formalis*), find its lexical base, retrieve all lemmas and affixes linked to that lexical base (word formation family)
- ► Find lemmas with a specific prefix (ad-) or suffix (-bil)
- ► List distinct lemmas with suffix -bil, prefix ad- ending in -bilis only

LiLa triplestore available at:
https://lila-erc.eu/data/

More sophisticated ways of querying the data:

- ► Start from a lemma (e.g. *formalis*), find its lexical base, retrieve all lemmas and affixes linked to that lexical base (word formation family)
- ► Find lemmas with a specific prefix (ad-) or suffix (-bil)
- ► List distinct lemmas with suffix -bil, prefix ad- ending in -bilis only
- ► Count the frequency of the 15 most used affixes attached to nouns.

Do you have research questions?
Contact me at eleonoramaria.litta@unicatt.it

Conclusions:

- ▶ WP models have a better way of demonstrating derivational relationships between words
- ▶ CxM was useful to develop a new model to account for Latin word formation that could fit into the LiLa Knowledge Base
- ▶ WFL in LiLa does not contain info on directionality.

Future plans:

- ▶ Find a way of defining and naming all "base" nodes
- ▶ More investigation is needed to add word formation specific semantic information to the LiLa knowledge base
- ▶ Enlarge the lexical basis for which WF is provided to Medieval Latin lemmas contained in Lemlat.

- ▶ Marco Pappalepore (for extracting triples from the old WFL)
- ▶ Francesco Mambrini (for creating the ontology behind the new WFL)

**Eleonora Litta**
Università Cattolica del Sacro Cuore, Milano

✉ eleonoramaria.litta@unicatt.it

○ https://github.com/CIRCSE/WFL

⊕ https://wfl.marginalia.it

📍 Largo Gemelli 1, 20123 Milan, Italy