

Paolo Monella

paolo.monella@gmx.net



In the Tower of Babel
Modelling primary sources of
multi-testimonial textual transmissions

Digital Classicist & Institute of Classical Studies Seminar 2012
London, 20 July 2012

Outline

- The 'Babel issue'
 - The 'Babel issue' in a nutshell
 - Comparing texts
 - In the Tower
- Two issues
 - A. Comparison at the linguistic layer
 - B. Comparison at the graphemic layer

Outline

- The 'Babel issue'
 - The 'Babel issue' in a nutshell
 - Comparing texts
 - In the Tower
- Two issues
 - A. Comparison at the linguistic layer
 - B. Comparison at the graphemic layer

The 'Babel issue' in a nutshell

- Each primary source uses a different writing (encoding, semiotic) system
- We want to compare the texts they bear

Outline

- The 'Babel issue'
 - The 'Babel issue' in a nutshell
 - Comparing texts
 - In the Tower
- Two issues
 - A. Comparison at the linguistic layer
 - B. Comparison at the graphemic layer

Comparing texts

- Each primary source uses a different writing (encoding, semiotic) system
- We want to compare the texts they bear
 - Textual Criticism
 - Processing (e. g. cross-corpus search)

Comparing texts

- Each primary source uses a different writing (encoding, semiotic) system
- We want to compare the texts they bear
 - Textual Criticism
 - Processing (e. g. cross-corpus search)

MS A
pui'_G

MS B
peruius_G

Print ed. C
pervius_G

Comparing texts

- Each primary source uses a different writing (encoding, semiotic) system
- We want to compare the texts they bear
 - Textual Criticism
 - Processing (e. g. cross-corpus search)

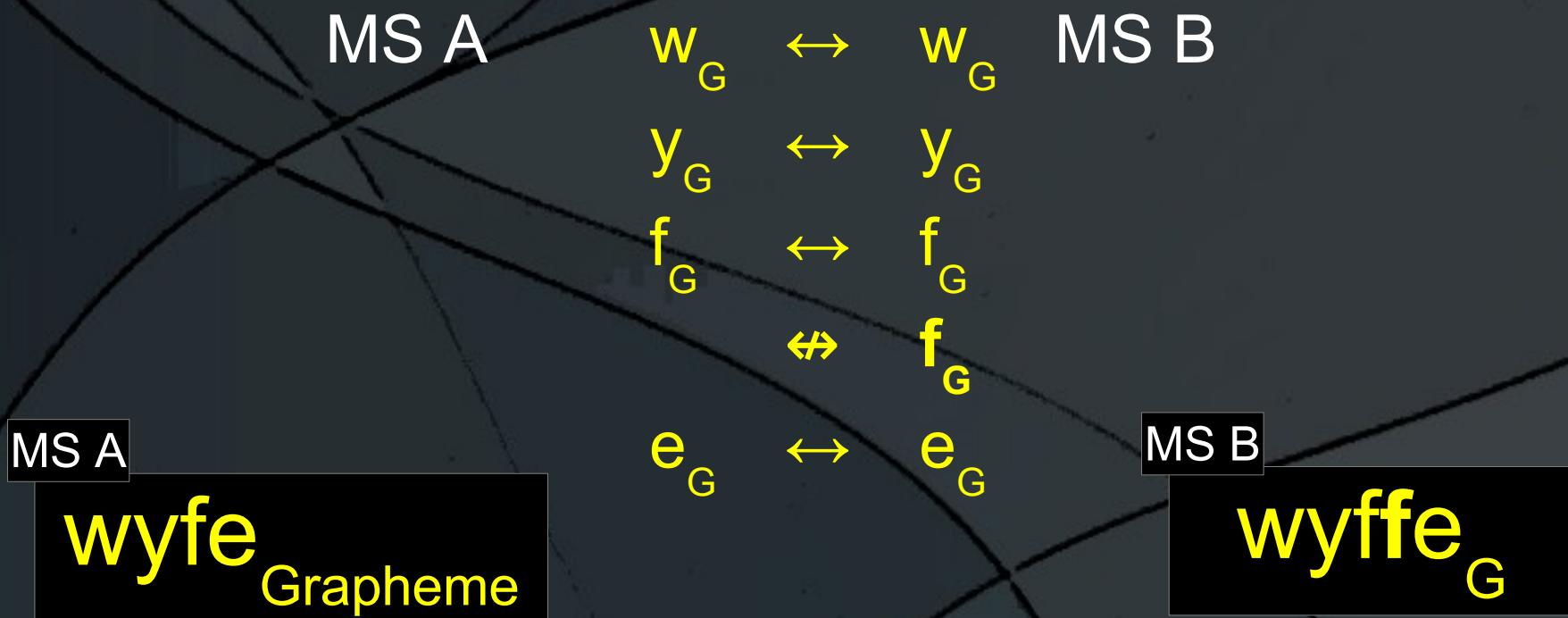
MS A
*pui'*_G

MS B
*peruius*_G

Query
*pervius*_G

Comparing texts

- Simplification: same writing system (no 'Babel' issue)
 - simple variant (at **graphemic** layer)



Comparing texts

- Simplification: same writing system (no 'Babel' issue)
 - no variant (at **linguistic** layer)

wyfe
Word

wyfe
Letter

wyffe
L

MS A

wyfe
Grapheme

MS B

wyffe
G

Comparing texts

- The same text (reading)?
 - Comparing texts at different layers

Comparing at
linguistic layer

wyfe_{Word}

wyfe_{Letter}

wyffe_L

MS A

wyfe_{Grapheme}

Comparing at
graphemic layer

wyffe_G

Comparing texts

- Types of edition in the print age

Critical edition (Pierazzo)
“Museum” edition (Vanhoutte)

wyfe_{Word}

wyfe_{Letter}

wyffe_L

MS A

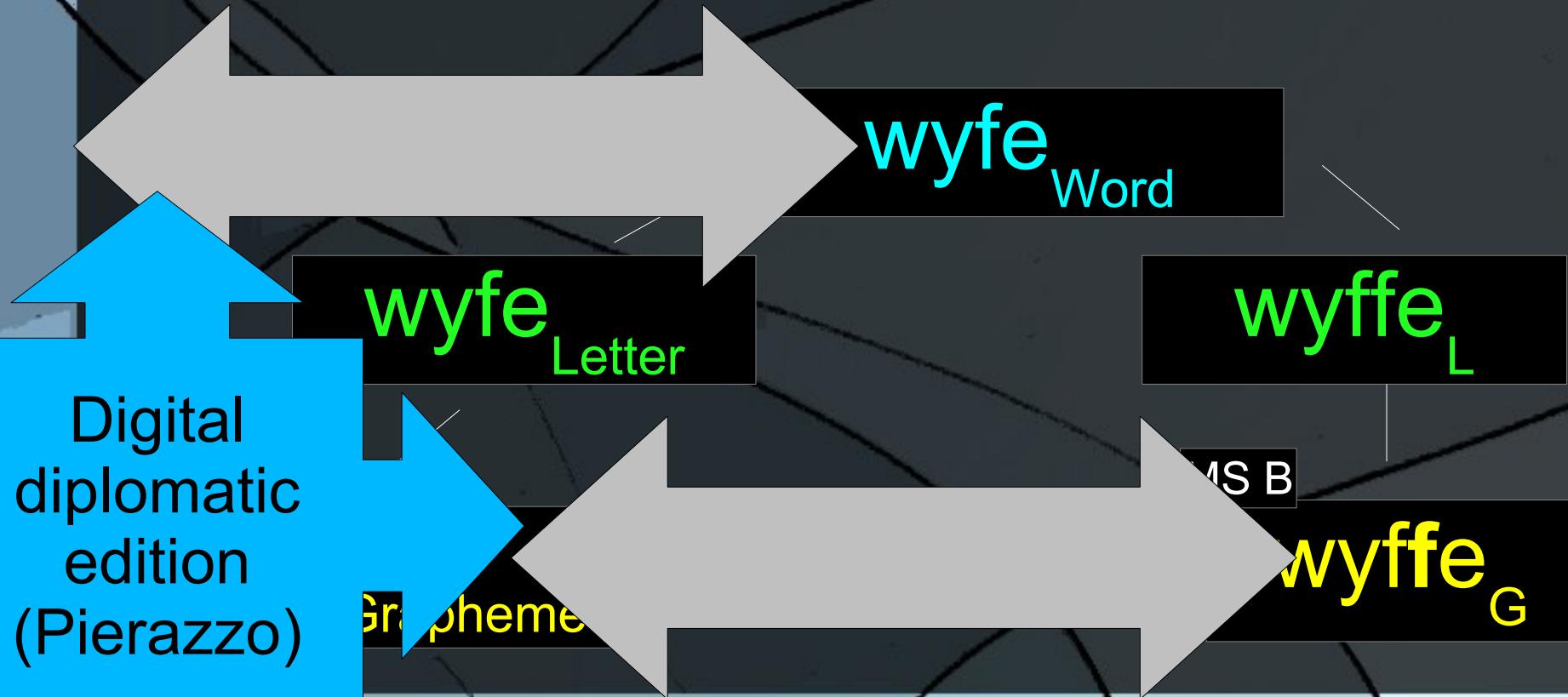
wyfe_{Grapheme}

Diplomatic edition (Pierazzo)
“Archive” edition (Vanhoutte)

wyffe_G

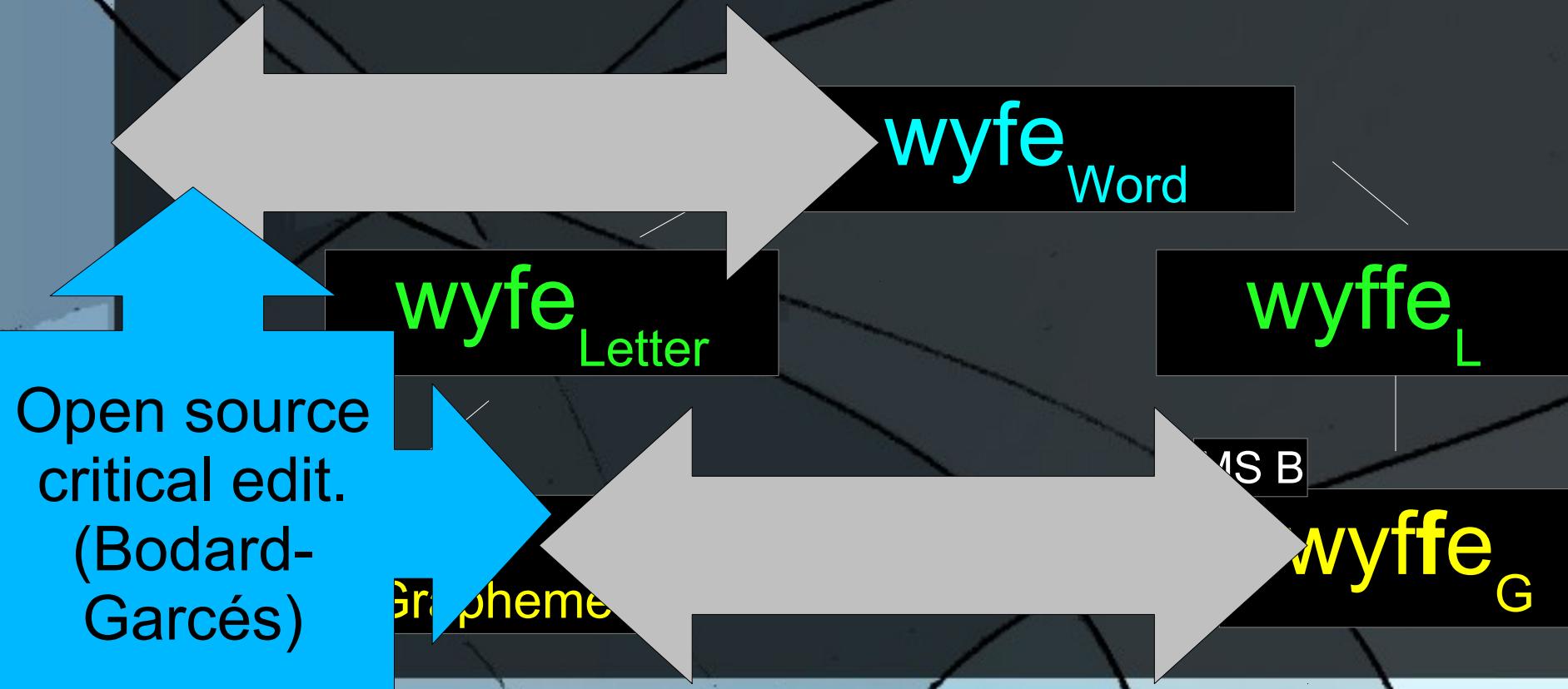
Comparing texts

- Digital edition



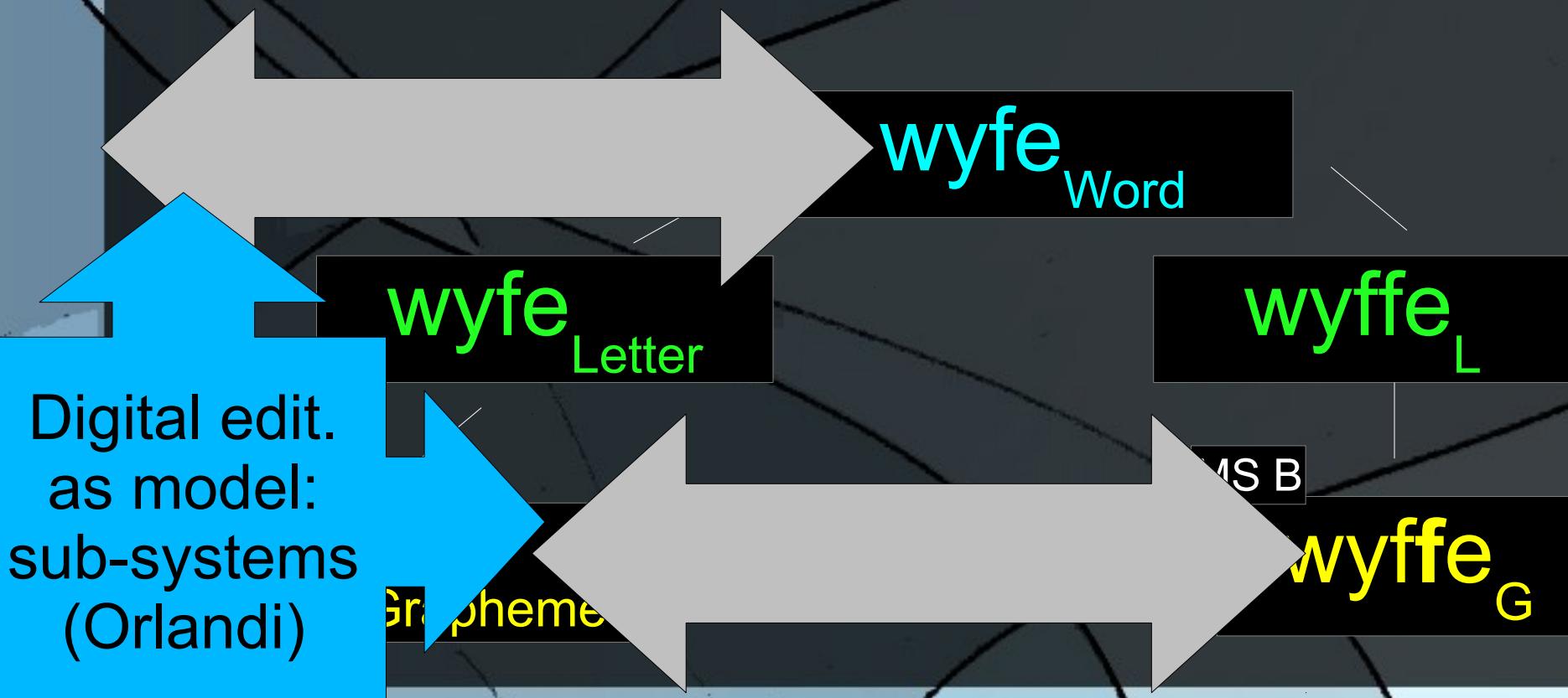
Comparing texts

- Digital edition



Comparing texts

- Digital edition



Outline

- The 'Babel issue'
 - The 'Babel issue' in a nutshell
 - Comparing texts
 - In the Tower
- Two issues
 - A. Comparison at the linguistic layer
 - B. Comparison at the graphemic layer

In the Tower

- Different Writing Systems (Graphemes):

A: p_G , u_G (no u_G/v_G distinction)

B: u_G (no u_G/v_G distinction)

C: u_G , v_G

peruius
Word

perUiUs
Letter

pervius
L

MS A

pui'
Grapheme

MS B

peruius
 G

PR C

pervius
 G

In the Tower

- Different Alphabets (**Letters**):

A: u_L (no u_L/v_L distinction)

B: u_L (no u_L/v_L distinction)

C: u_L, v_L

peruius
Word

perUiUs
Letter

pervius_L

MS A

pui'
Grapheme

MS B

peruius_G

PR C

pervius_G

In the Tower

- Same text at Linguistic layer (**Inflected word**)
“Perius”, nom. sing. masc. of lemma “perius,
-a, um”

peruius
Word

perUiUs
Letter

pervius
L

pui'
Grapheme

peruius
G

pervius
G

MS A

MS B

PR C

In the Tower

- Digital edition: Formalisation

(some code)_C

2 ..._C

2 ..._C

MS A

&#A751; K E
&us;_C

MS B

2 - 4
K E
K I_C

PR C

2 - 4
L E
K I_C

In the Tower

- Life in the Tower of Babel (recapitulation):
 - Each builder, a language
 - Each primary source, a semiotic system
 - ...yet we want builders to interact
 - Textual Criticism
 - Processing (e. g. cross-corpus search)
 - ...at different layers
 - Graphs, allographs, graphemes, linguistic...
 - ... all the sudden, all builders are robots!
 - Formalisation (no human intuition)

Outline

- The 'Babel issue'
 - The 'Babel issue' in a nutshell
 - Comparing texts
 - In the Tower
- Two issues
 - A. Comparison at the linguistic layer
 - B. Comparison at the graphemic layer

Two issues

A. Comparison at **linguistic** layer

- Substantial readings ($\text{wyfe}_W \neq \text{lyf}_W$)
- Critical edition
- The 'text' (**Inflected words**)

B. Comparison at **graphemic** layer

- Accidentals ($\text{wyfe}_G \neq \text{wyffe}_G$)
- Diplomatic edition
- The 'spelling' (orthography, **Graphemes**)

Two issues

A. Comparison at **linguistic** layer

- How do you identify elements at linguistic layer (**Inflected words**) **digitally**?
 - The Canterbury Tales Project:
“Regularized spelling”
 - Tito Orlandi:
Linguistic entities

Two issues

B. Comparison at graphemic layer

- Can you compare **graphemes** through different graphemic systems?
 - The Canterbury Tales Project:
Unicode; Corpus-wide modelling of the graphemic system
 - Tito Orlandi:
MS-wide complete modelling of the graphemic system (and of other systems at other textual layers)

Two issues

- The Canterbury Tales Project
- Tito Orlandi, *Informatica Testuale*, Laterza:
Roma 2010
- My own project of an experimental edition
from the *Anthologia Latina*

Outline

- The 'Babel issue'
 - The 'Babel issue' in a nutshell
 - Comparing texts
 - In the Tower
- Two issues
 - A. Comparison at the linguistic layer
 - B. Comparison at the graphemic layer

A. Comparison at linguistic layer

- Why comparing Words (at linguistic layer)?
The Canterbury Tales Project:
 - Textual criticism
 - Relations between MSS
 - Processing
 - Indexing (“Spelling database”)



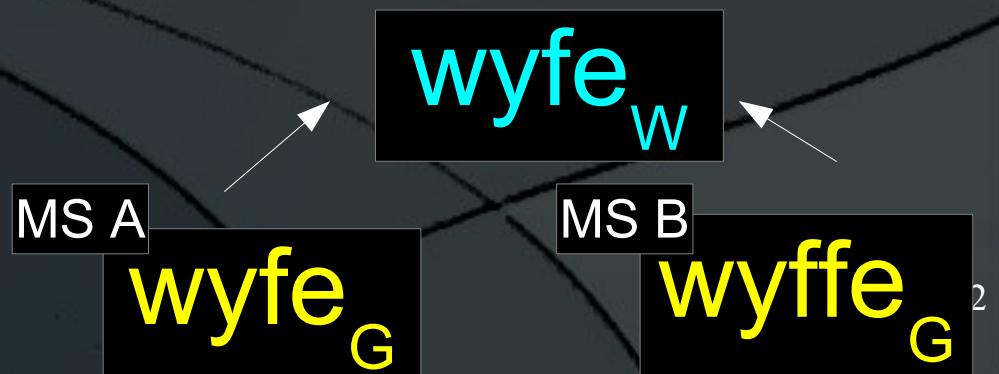
A. Comparison at linguistic layer

- From **Grapheme** to **Word**, not **Lemma**
 - Inflected word ($wyfe_w \neq wyves_w$)
 - Not as element of language (*langue*)
 - But as as element of that text (*parole*)
 - “Substantial reading”
 - Critical edition



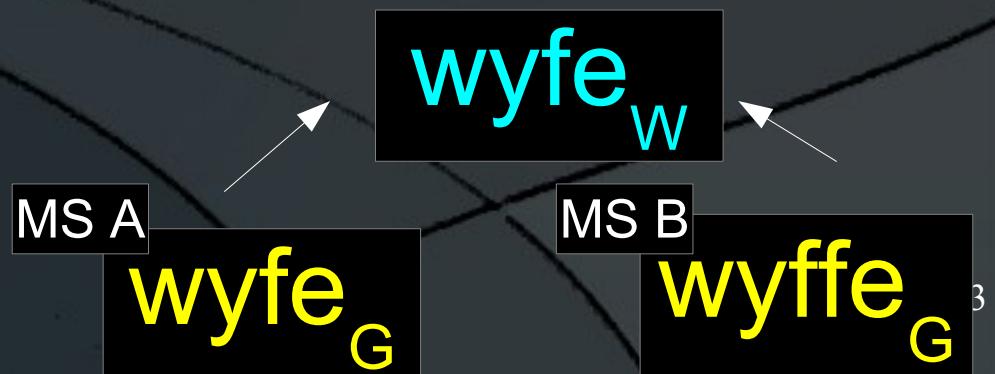
A. Comparison at linguistic layer

- From Grapheme to Word: semi-automatic
 - Not computable so far
 - Linguistic competence required
 - Grapheme → Letter: ambiguous graphemes
 $\text{lon}_G \rightarrow \text{Jon}_L / \text{Ion}_L \rightarrow \text{Jon}_W / \text{Ion}_W$
 $\text{pdo}_G \rightarrow \text{perdo}_L / \text{proto}_L \rightarrow \text{perdo}_W / \text{proto}_W$



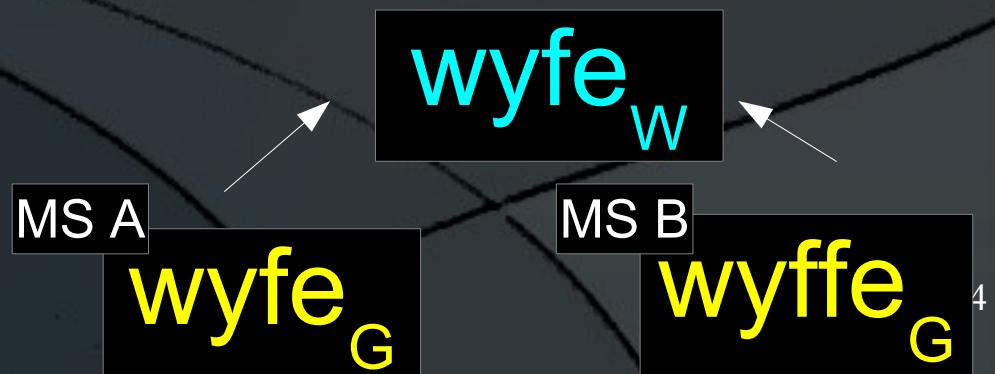
A. Comparison at linguistic layer

- From Grapheme to Word: semi-automatic
 - Not computable so far
 - Linguistic competence required
 - Letter → Word: omographs
 $\text{est}_G \rightarrow \text{est}_L \rightarrow \text{est} (\text{she is})_W / \text{est} (\text{she eats})_W$



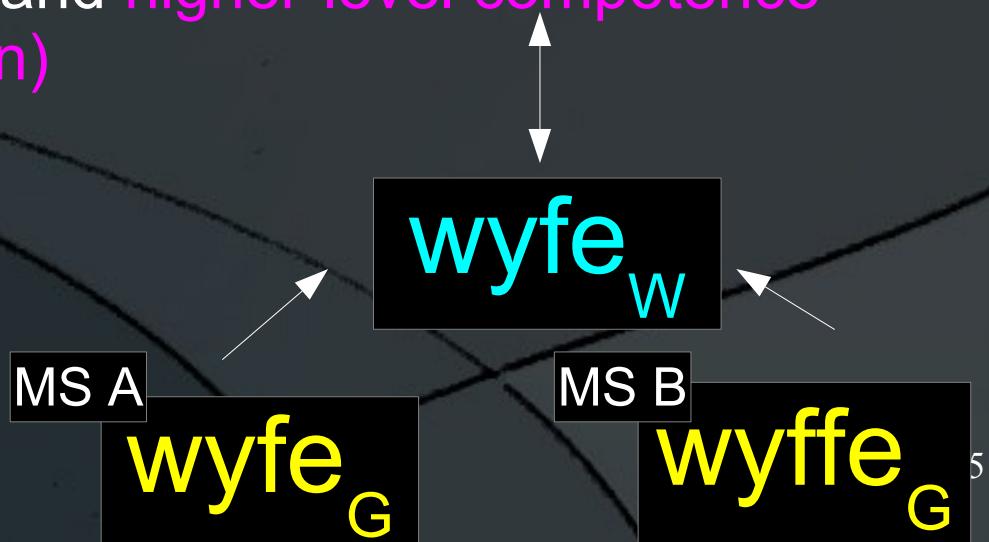
A. Comparison at linguistic layer

- From Grapheme to Word: semi-automatic
 - Not computable so far
 - Linguistic competence required
 - Letter → Word: different spellings
 $wyfe_w / wyffe_w \rightarrow wyfe_w / wyffe_w \rightarrow wyfe_w$



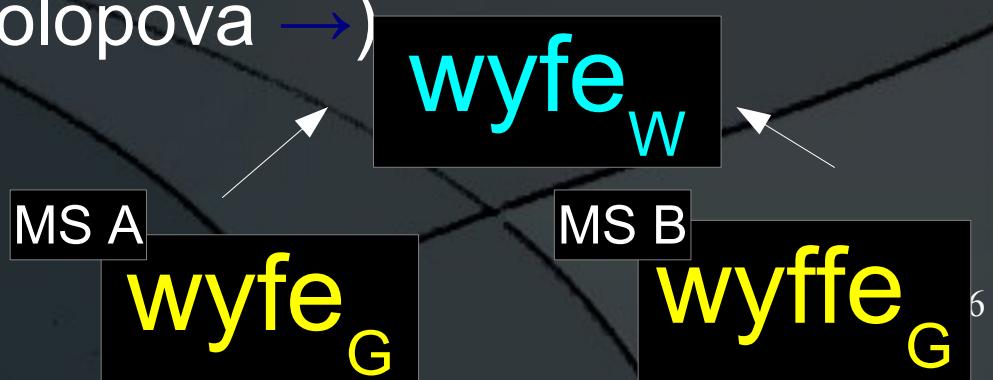
A. Comparison at linguistic layer

- From Grapheme to Word: semi-automatic
 - Not computable so far
 - Linguistic competence required
 - 'Context': interpretation of the whole text
 - Linguistic and higher-level competence (Jon / Ion)



A. Comparison at linguistic layer

- From Grapheme to Word: semi-automatic
 - Semi-automatic procedures (human-driven, computer-assisted)
 - "The computer collation program we are using (Collate) permits regularization as part of the collation process" [...] "regularized spelling" (Robinson-Solopova →)



A. Comparison at linguistic layer

- How do you identify elements at linguistic layer (**Inflected words**) digitally?

- The Canterbury Tales Project:
“Regularized spelling”

1000101_c

- Tito Orlandi:
Linguistic entities

0010100_c

MS A

wyfe_G

MS B

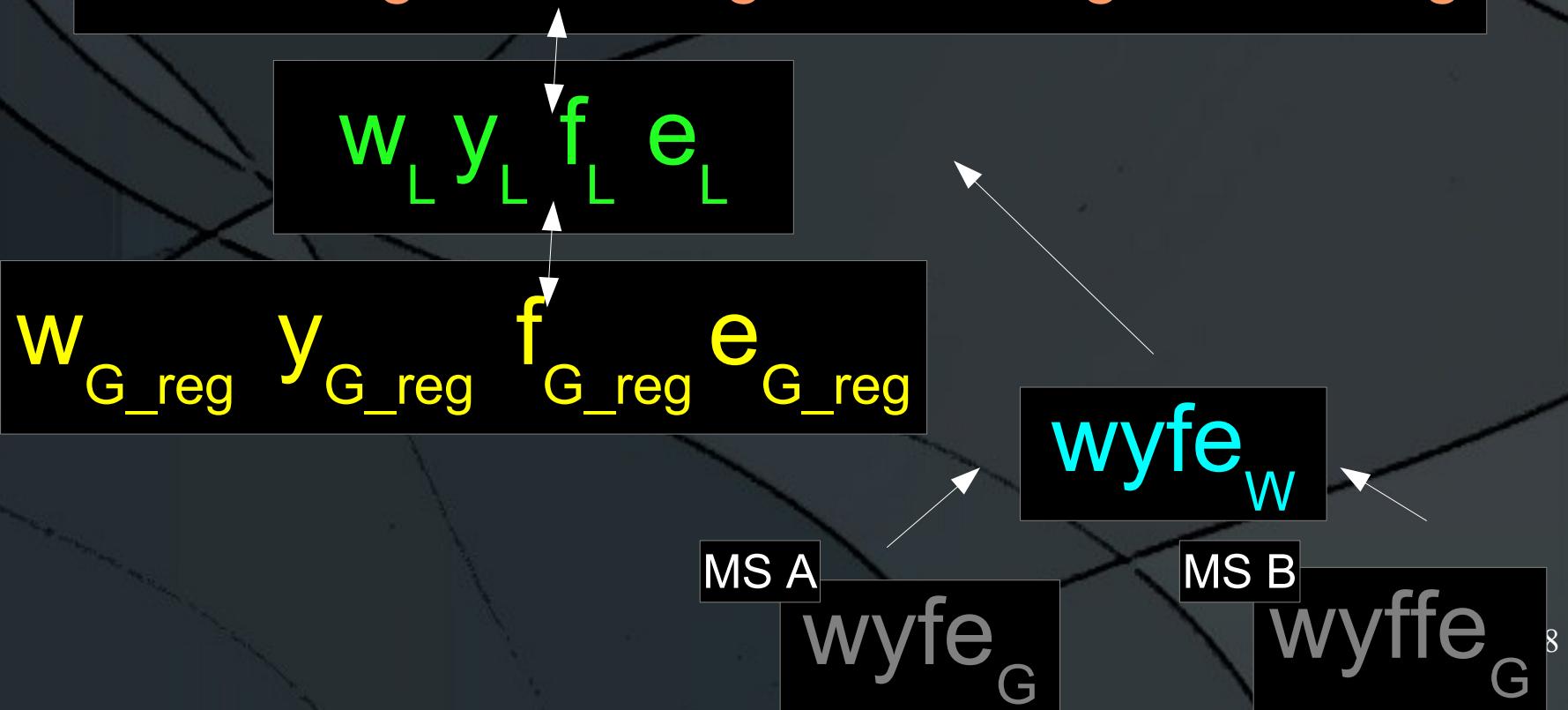
wyffe_G⁷

wyfe_W

A. Comparison at linguistic layer

- Canterbury: “Regularized spelling”

U0077_c U0079_c U0066_c U0065_c



A. Comparison at linguistic layer

- Orlandi: Discrete linguistic entities in a database

ID_09834_C

wyfe_W

MS A

wyfe_G

MS B

wyffe_{G⁹}

A. Comparison at linguistic layer

- Orlandi: Discrete linguistic entities in a database → How do you identify them?

Lemma [wyfe]
Morph [sing]_C



A. Comparison at linguistic layer

- Orlandi: Discrete linguistic entities in a database → How do you identify them?

[deus1]
[gen.pl.masc]_C

[deus1]
[gen.pl.masc]_C

deorum_W

MS A

deorum_G

deum_W

MS B

devm_G

A. Comparison at linguistic layer

- Orlandi: Discrete linguistic entities in a database → How do you identify them?

[deus1]
[gen.pl.masc.1]_C

[deus1]
[gen.pl.masc.2]_C

deorum_W

MS A

deorum_G

deum_W

MS B

devm_G

A. Comparison at linguistic layer

- Orlandi: Discrete linguistic entities in a database → How do you identify them?

deorum_{G_reg} → U0064_c U0065_c
U006F_c U0072_c U0075_c U006D_c

[deus1]
[gen.pl.masc]_c

deorum_w

MS A

deorum_G

deum_{G_reg} → U0064_c U0065_c
U0075_c U006D_c

[deus1]
[gen.pl.masc]_c

deum_w

MS B

devm_G

A. Comparison at linguistic layer

- Orlandi/Monella: is it worth it?

Canterbury

deum_{G_reg}

Orlandi/Monella

deum_{G_reg}
[deus1] [gen.plur.masc]_C

deum_W

devm_G

A. Comparison at linguistic layer

Canterbury, MS A

deum_{G_reg}

Canterbury, MS B

deum_{G_reg}

Orlandi/Monella, MS A

deum_{G_reg}
[deus1] [gen.plur.masc]_C

Orlandi/Monella, MS B

deum_{G_reg}
[deus1] [acc.sing.masc]_C

MS A] *Pater devm*

MS B] *Timeo devm*

deum_W

≠ deum_W

A. Comparison at linguistic layer

- It's only worth it
 - if words *are not* their regularised spelling, i. e.
 - if homograph words exist

deum_w

Pater deum

≠ deum_w

Timeo deum

A. Comparison at linguistic layer

- It's only worth it
 - if substantial readings are *not* their regularised spelling, i. e.
 - if homograph substantial readings exist (do they? → Concept of “reading”)

MS A

*Pater,
devm Neptvnvm odi*

deum

MS B

*Pater deum,
Neptunum odi*

deum?

A. Comparison at linguistic layer

Canterbury, MS A

deum_{G_reg}

Canterbury, MS B

deum_{G_reg}

Orlandi/Monella, MS A

deum_{G_reg}
[deus1] [gen.plur.masc]_C

Orlandi/Monella, MS B

deum_{G_reg}
[deus1] [acc.sing.masc]_C

MS A

*Pater,
devm Neptvnvm odi*

MS B

*Pater deum,
Neptunum odi*

deum

≠ deum?

A. Comparison at linguistic layer

Canterbury, MS A

deum_{G_reg}

Canterbury, MS B

deum_{G_reg}

Orlandi/Monella, MS A

deum_{G_reg}
[deus1] [gen.plur.masc]_C

Orlandi/Monella, MS B

deum_{G_reg}
[deus1] [acc.sing.masc]_C

MS A

Pater,
devm Neptvnvm odi

MS B

Pater deum,
Neptunum odi

deum

= deum?

A. Comparison at linguistic layer

- Normally done via regularised spelling
- Inflected words are not their regularised spelling (omographs)
- Are there omograph substantial readings (*deum*)?
- If so:
 - No regularised spelling
 - But formal identifiers for inflected words (spelling, lemma, identifier)

Outline

- The 'Babel issue'
 - The 'Babel issue' in a nutshell
 - Comparing texts
 - In the Tower
- Two issues
 - A. Comparison at the linguistic layer
 - B. Comparison at the graphemic layer

B. Comparison at graphemic layer

- Why comparing spellings (at graphemic layer)?
 - Historical Linguistics
 - Evidence for lexical, morphological and phonetic evolution

B. Comparison at graphemic layer

- Why comparing spellings (at graphemic layer)?
 - Palaeography
 - Indirect evidence for letters (Benskin: letter þ *thorn* →)
 - Overlapping with Historical Linguistics (Emiliano: “Scripto-linguistic change” →)
 - Though only graphemes
 - Not allographs (graphetes, s/l) or graphs (bitmap)

B. Comparison at graphemic layer

- Why comparing spellings (at graphemic layer)?
 - Textual Criticism
 - 'Orthographic' apparatus (The Hengwrt Chaucer Digital Facsimile, Collation function →)
 - “Although for most manuscripts collation of the regularized text will produce sufficient information **to place those manuscripts in genetic relation to one another [...]**” (Robinson-Solopova →)

B. Comparison at graphemic layer

- *The loue of love*
- Loue: “A hill or mountain” ≠ Loye: “Love”
- Middle English alphabet: u ≠ v



B. Comparison at graphemic layer

a _L
u _L
v _L

The loue of love

MS A

The loue of loue

MS B⁵⁶

B. Comparison at graphemic layer

a_L
u_L
v_L

The loue of love_G

MS A

The loue of loue_G

MS B ⁵⁷

B. Comparison at graphemic layer

a_G	a_L	a_G
u_G	u_L	u_G
v_G	v_L	

The $lou_G e$ of $lov_G e$
MS A

The $lou_G e$ of $lou_G e$
MS B

B. Comparison at graphemic layer

a_G	a_L	a_G
u_G	u_L	u_G
v_G	v_L	

- Back in the Tower



The lou_Ge of lov_Ge
MS A

The lou_Ge of lou_Ge
MS B 59

B. Comparison at graphemic layer

a_G	a_L	a_G
u_G	u_L	u_G
v_G	v_L	

- Back in the Tower
 - “Grapheme. (1) A minimally distinctive unit of writing in the context of a particular writing system” (Unicode Glossary →).

The $lou_G e$ of $lov_G e$

MS A

The $lou_G e$ of $lou_G e$

MS B⁶⁰

B. Comparison at graphemic layer

a_G	a_L	a_G
u_G	u_L	u_G
v_G	v_L	

- Back in the Tower
 - De Saussure: relational nature of signs *within* a semiotic system

The $lou_G e$ of $lov_G e$
MS A

The $lou_G e$ of $lou_G e$
MS B⁶¹

B. Comparison at graphemic layer

a_G	a_L	a_G
u_G	u_L	u_G
v_G	v_L	

- Back in the Tower
 - De Saussure: relational nature of signs *within* a semiotic system
 - This is a different grapheme than this, as the former is in contrast with this, while the latter is not

The $\text{lou}_G e$ of $\text{lov}_G e$

MS A

The $\text{lou}_G e$ of $\text{lou}_G e$

MS B ⁶²

B. Comparison at graphemic layer

a_G	a_L	a_G
u_G	u_L	u_G
v_G	v_L	

- Back in the Tower
 - De Saussure: relational nature of signs *within* a semiotic system
 - This is a different grapheme than this, as the former is in contrast with this, while the latter is not

The $lou_G e$ of $lov_G e$

MS A

The $lou_G e$ of $lou_G e$

MS B 63

B. Comparison at graphemic layer

a_G	a_L	a_G
u_G	u_L	u_G
v_G	v_L	

- Back in the Tower
 - De Saussure: relational nature of signs *within* a semiotic system
 - This is a different grapheme than this, as the former is in contrast with this, while the latter is not

The $lou_G e$ of $lov_G e$

MS A

The $lou_G e$ of $lou_G e$

MS B⁶⁴

B. Comparison at graphemic layer

a_G	a_L	a_G
u_G	u_L	u_G
v_G	v_L	

- How do you encode graphemes digitally?

The lou_Ge of lov_Ge
MS A

The lou_Ge of lou_Ge⁶⁵
MS B

B. Comparison at graphemic layer

a_G	a_L	a_G
u_G	u_L	
v_G	v_L	u_G

- TEI
 - Unicode → →
- The Canterbury Tales Project
 - Corpus-wide definition of the graphemic system →

The lou_Ge of lov_Ge
MS A

The lou_Ge of lou_Ge
MS B⁶⁶

B. Comparison at graphemic layer

a_G	a_L	a_G
u_G	u_L	
v_G	v_L	u_G

U0075_c

U0076_c

The lou_Ge of lov_Ge

MS A

U0075_c

U0075_c

The lou_Ge of lou_Ge

MS B⁶⁷

B. Comparison at graphemic layer

a_G	a_L	a_G
u_G	u_L	
v_G	v_L	u_G

No match:
different graphemes

U0075_c

U0076_c

The lou_Ge of lov_Ge

MS A

U0075_c

U0075_c

The lou_Ge of lou_Ge

MS B⁶⁸

B. Comparison at graphemic layer

a_G	a_L	a_G
u_G	u_L	
v_G	v_L	u_G

Match:
same grapheme

U0075_c

U0076_c

The lou_Ge of lov_Ge

MS A

U0075_c

U0075_c

The lou_Ge of lou_Ge

MS B⁶⁹

B. Comparison at graphemic layer

a_G	a_L	a_G
u_G	u_L	u_G
v_G	v_L	

- Orlandi
 - MS-wide *complete* definition of graphemic system
 - Not corpus-wide (Canterbury)
 - Not world-wide (Unicode)

The lou_Ge of lov_Ge
MS A

The lou_Ge of lou_Ge
MS B⁷⁰

B. Comparison at graphemic layer

a_G	a_L	a_G
u_G	u_L	u_G
v_G	v_L	

<teiHeader>

<encodingDesc>

<charDecl>

Graphemes

→ <char xml:id="uv">

Allographs

→ <glyph xml:id="long_s">

The lou_Ge of lov_Ge

MS A

The lou_Ge of lou_Ge

MS B⁷¹

B. Comparison at graphemic layer

a_G	a_L	a_G
u_G	u_L	
v_G	v_L	u_G

<body>

<p>

<g ref="uv" />

<!-- or, better: -->

<!ENTITY uv '<g ref="#uv" />

&uv;

The lou_Ge of lov_Ge

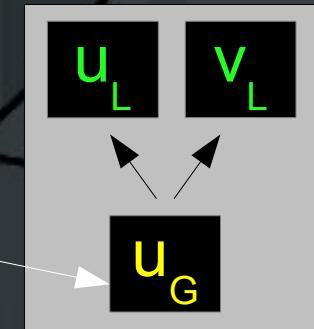
MS A

The lou_Ge of lou_Ge

MS B

B. Comparison at graphemic layer

a_G	a_L	a_G
u_G	u_L	u_G
v_G	v_L	

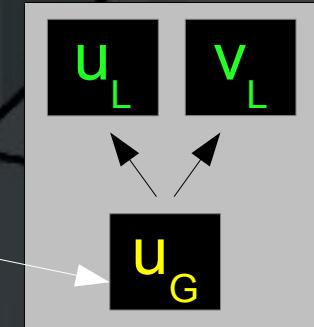


The $lou_G e$ of $lov_G e$
MS A

The $lou_G e$ of $lou_G e$
MS B 73

B. Comparison at graphemic layer

a_G	a_L	a_G
u_G	u_L	u_G
v_G	v_L	



MS A

```
<charDecl>
  <char xml:id= "a">
  <char xml:id= "u">
  <char xml:id= "v">
```

MS B

```
<charDecl>
  <char xml:id= "a">
  <char xml:id= "uv">
```

U0075_C

U0076_C

The lou_Ge of lov_Ge

MS A

&uv_C

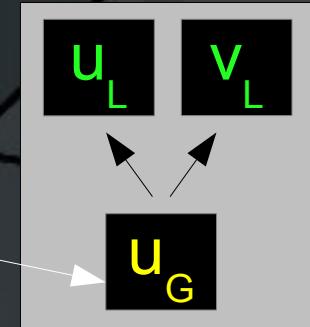
&uv_C

The lou_Ge of lou_Ge

MS B⁷⁴

B. Comparison at graphemic layer

a_G	a_L	a_G
u_G	u_L	u_G
v_G	v_L	



MS B

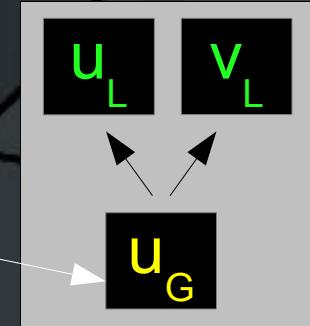
```
<charDecl>
  <char xml:id= "uv">
    <charName>SMALL LATIN U OR V</charName>
    <desc>U-shaped when lowercase, V-shaped when
      uppercase. Content: either small letter
      Latin u or small letter Latin v</desc>
```

The lou_Ge of lov_Ge
MS A

The lou_Ge of lou_Ge
MS B

B. Comparison at graphemic layer

a_G	a_L	a_G
u_G	u_L	u_G
v_G	v_L	



MS B

```
<charDecl>
<char xml:id= "uv">
<charProp>
  <localName>Expression</localName>
  <value>U+0075</value>
  <localName>Content</localName>
  <value>u | v</value>
```

The lou_Ge of lov_Ge

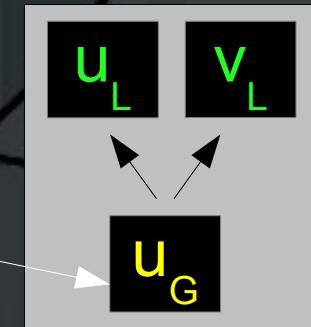
MS A

The lou_Ge of lou_Ge

MS B⁷⁶

B. Comparison at graphemic layer

a_G	a_L	a_G
u_G	u_L	u_G
v_G	v_L	



The $lou_G e$ of $lov_G e$
MS A

The $lou_G e$ of $lou_G e$
MS B ⁷⁷

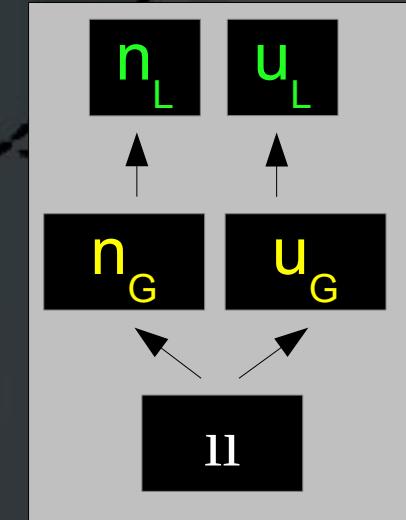
B. Comparison at graphemic layer

Two minims, lowercase

Expression (shape): indistinguishable.
Content (letter): n or u.

The reader does not identify the grapheme from its shape, but guessing its content.

Graphemic information is not conveyed by graphic information, but by linguistic information (context), so the scribe was confident in our Linguistic competence to tell the graphemes apart.



The lou_Ge of lov_Ge

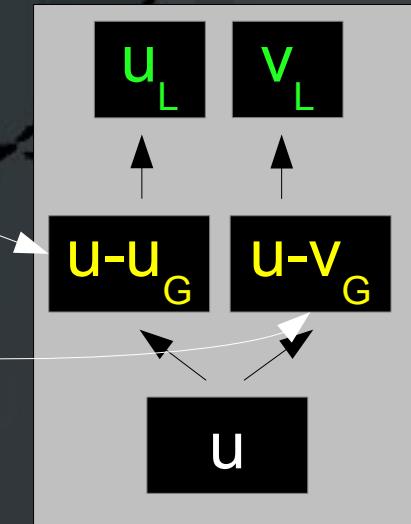
MS A

The lou_Ge of lou_Ge

MS B

B. Comparison at graphemic layer

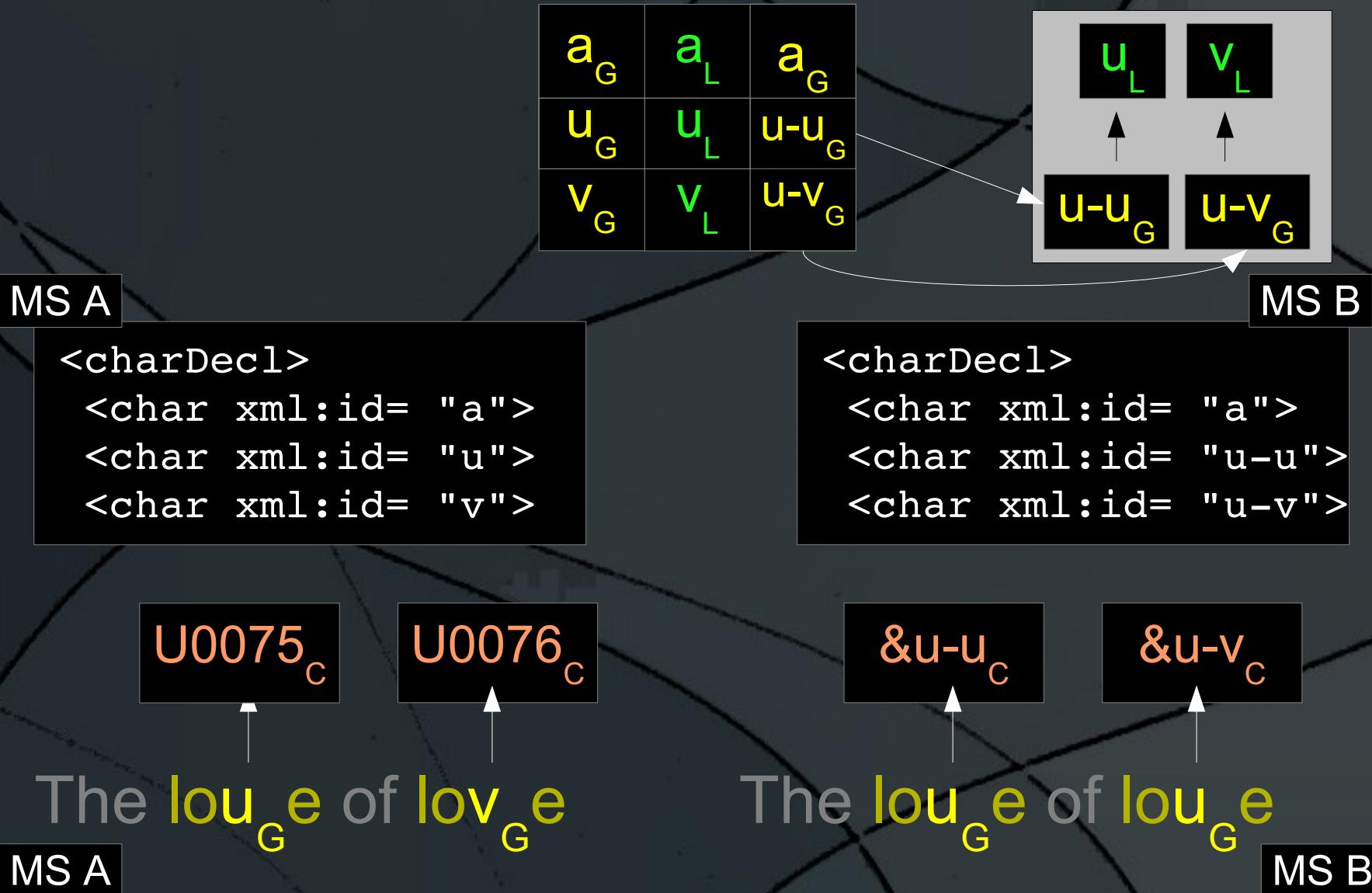
a_G	a_L	a_G
u_G	u_L	$u-u_G$
v_G	v_L	$u-v_G$



The $\text{lo}_G \text{e}$ of $\text{lov}_G \text{e}$
MS A

The $\text{lo}_G \text{e}$ of $\text{lou}_G \text{e}$
MS B ⁷⁹

B. Comparison at graphemic layer



B. Comparison at graphemic layer

a_G	a_L	a_G
u_G	u_L	$u-u_G$
v_G	v_L	$u-v_G$

- Fun
- But how do you compare graphemes *now*?

U0075_c

U0076_c

The lou_Ge of lov_Ge

MS A

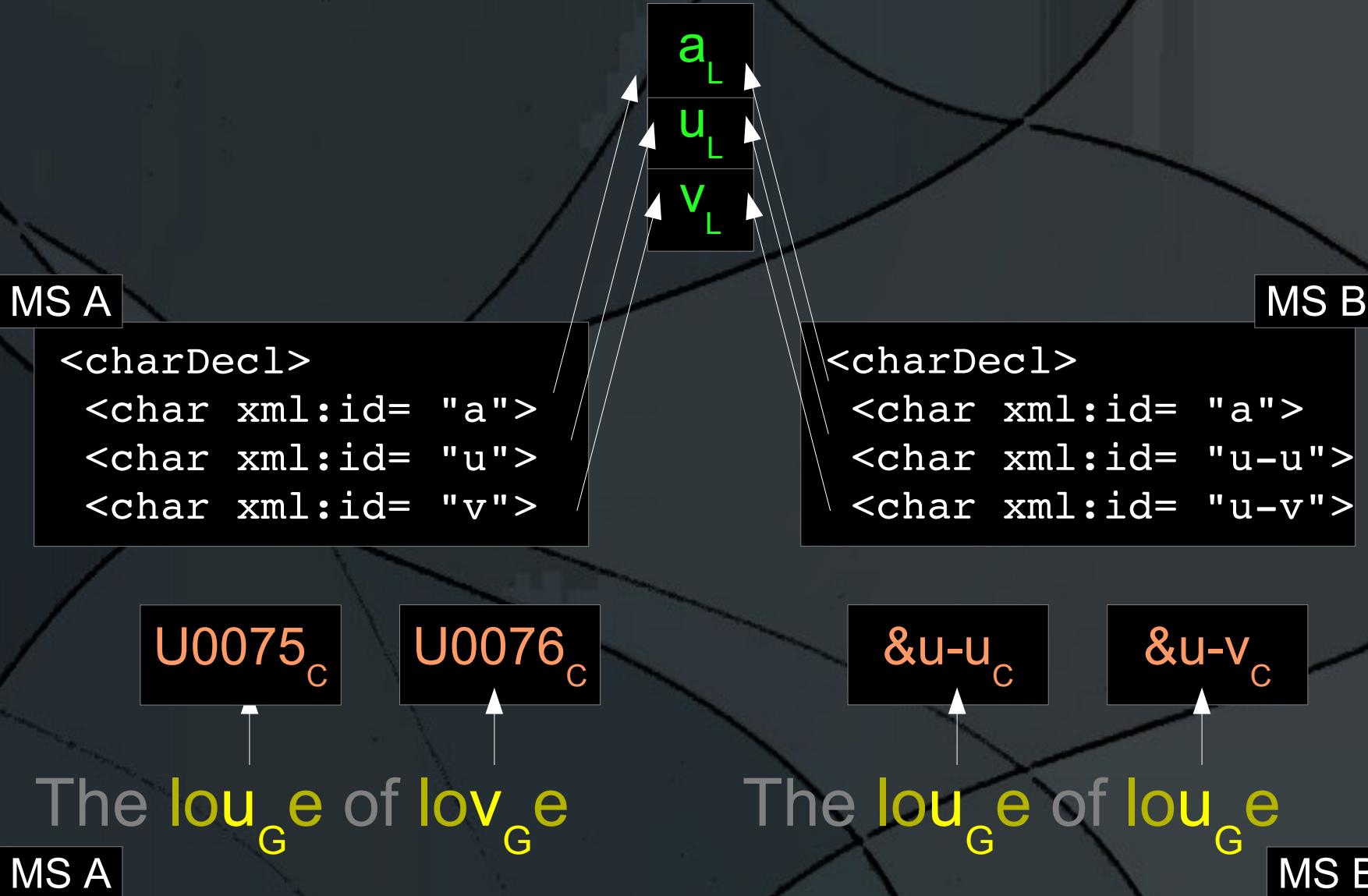
&u-u_c

&u-v_c

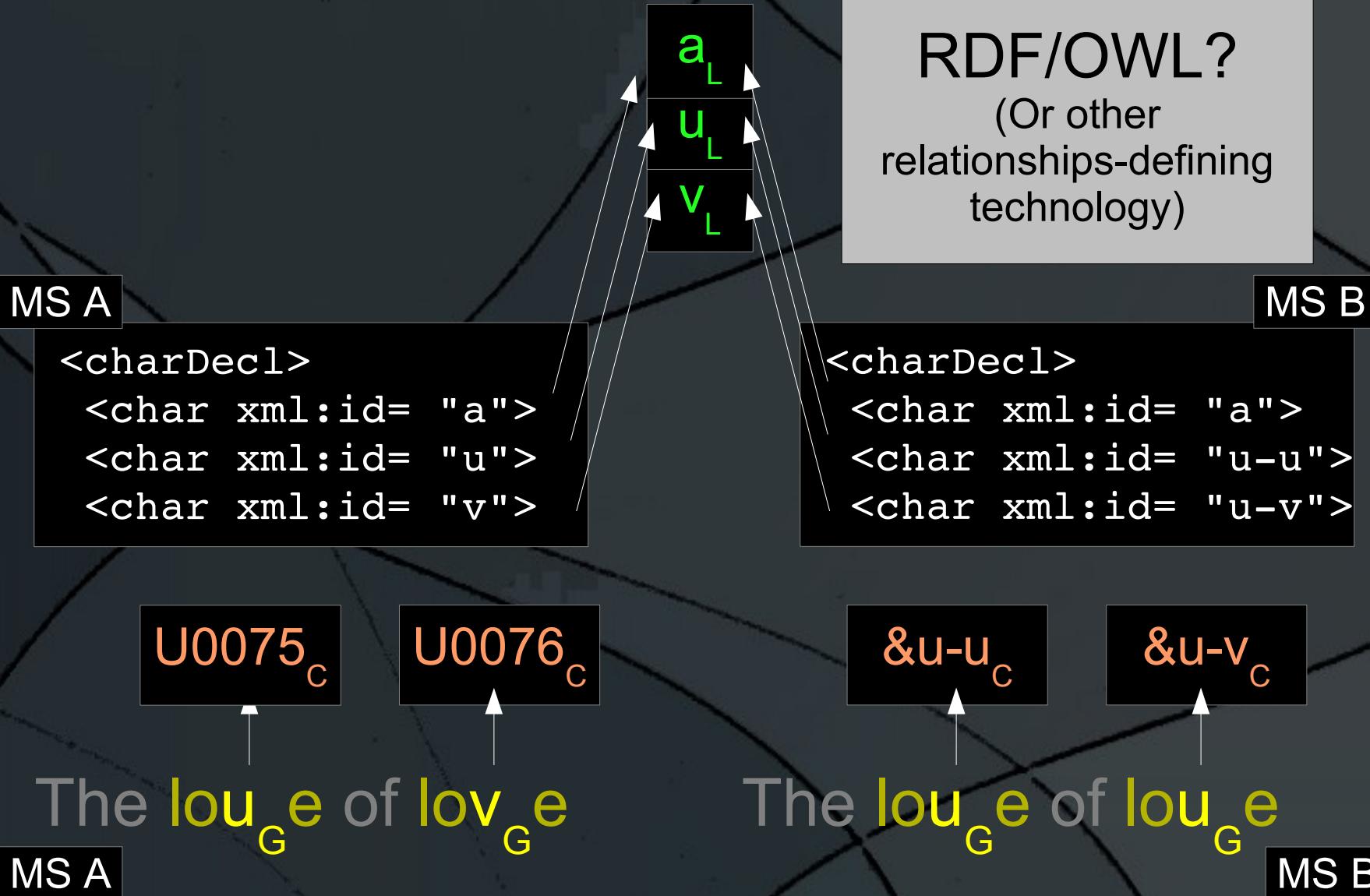
The lou_Ge of lou_Ge

MS B

B. Comparison at graphemic layer



B. Comparison at graphemic layer



B. Comparison at graphemic layer

a_G	a_L
h_G	h_L
t_G	t_L

a_L	a_G
h_L	h_G
t_L	t_G
b_L	b_G

- When the game gets tough...
 - Different alphabets
(Middle English, thorn letter)

B. Comparison at graphemic layer



- When the game gets tough...
 - Different alphabets, shared phonetics
(Middle English, thorn letter)

B. Comparison at graphemic layer

a_G	a_L	a_G
p_G	p_L	p_G
e_G	e_L	e_G
r_G	r_L	r_G
		p_G

- When the game gets tougher...
 - No 1:1 grapheme / letter ratio (Huitfeldt →)
 - Brevigraphs as graphemes (up to 40%)
 - “Grapheme. (1) A minimally distinctive unit of writing in the context of a particular writing system” (Unicode Glossary →).

B. Comparison at graphemic layer

a_G	a_L	a_G
p_G	p_L	p_G
e_G	e_L	e_G
r_G	r_L	r_G
		p_G

- When the game gets tougher...
 - No 1:1 grapheme / letter ratio (Huitfeldt →)
 - Brevigraphs as graphemes (up to 40%)
 - “Grapheme. (1) A minimally distinctive unit of writing in the context of a particular writing system” (Unicode Glossary →).

B. Comparison at graphemic layer

a_G	a_L
u_G	u_L
v_G	v_L

a_P
u_P
w_P

a_P
u_P
w_P

a_L	a_G
v_L	v_G

- When the game gets most tough...
 - (Classicalists start to play)
 - Different alphabets, different phonetics
(Ancient Latin)

B. Comparison at graphemic layer

a_G	a_L
u_G	u_L
v_G	v_L

a_L	a_G
v_L	v_G

- When the game gets most tough...
 - (Classicalists start to play)
 - Different alphabets, different phonetics
(Ancient Latin)
 - Still possible to formalise relationships

B. Comparison at graphemic layer

- Description of graphemic system
 - MS-wide
 - Complete
- Formalisation of relationships between graphemes in different graphemic systems
 - Involving higher-level entities (letters, phonemes)
- Intelligent searches, intelligent results

Outline

- The 'Babel issue'
 - The 'Babel issue' in a nutshell
 - Comparing texts
 - In the Tower
- Two issues
 - A. Comparison at the linguistic layer
 - B. Comparison at the graphemic layer