Paolo Monella

# Encoding Pre-modern Writing Systems

*ERC PAGES (AdG 2019 n° 882588)*

# Outline

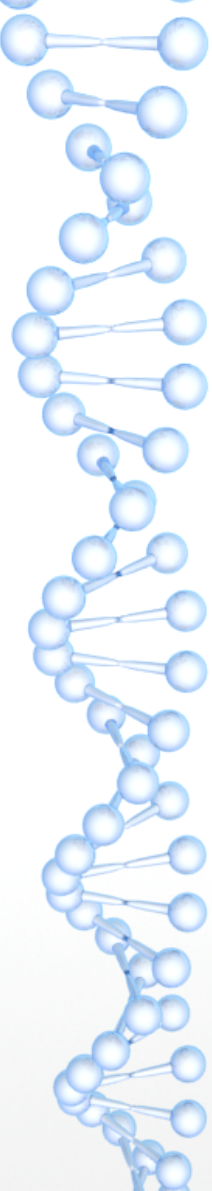- **Interoperability**
  of digital scholarly editions (DSEs)
  based on diplomatic transcriptions
  - The issue
  - Current solutions
  - Interoperability through modelling
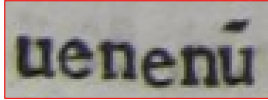- Orlandi's **table of sign**
- **Graphemes/allographs**

# Interoperability: the issue

# Interoperability: the issue

uenenú

- uenenū

# Interoperability: the issue

uenenú

- uenenū

**Diplomatic**

- Manual or HTR
- Visualization
- Processing

# Interoperability: the issue

uenenú

- uenenū

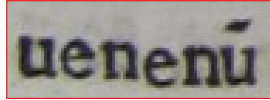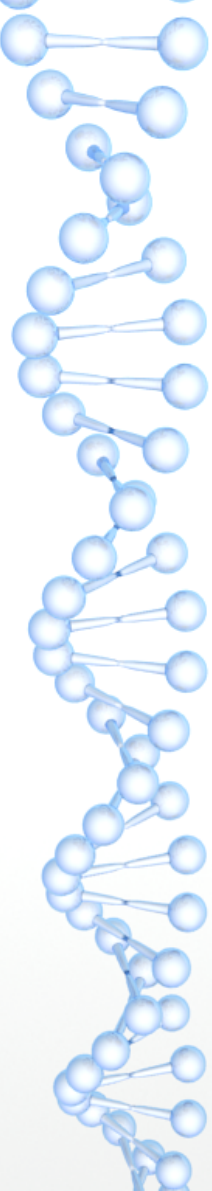# Interoperability: the issue



- uenenū

- uenenum

# Interoperability: the issue

- Processing
  - Search
  - Collation
  - NLP (lemma, PoS etc.)
  - Statistics (dist. reading)

uenenú

- uenenū

- uenenum

# Interoperability: the issue

- Processing
  - **Search**
  - Collation
  - NLP (lemma, PoS etc.)
  - Statistics (dist. reading)

uenenū

- uenenū

- uenenum

venenum

# Interoperability: the issue



- Processing
  - Search
  - **Collation**
  - NLP (lemma, PoS etc.)
  - Statistics (dist. reading)

uenenū

- uenenū

- uenenum

112 excedit] excedit corr. ex
exceditis R n70
114 obicitur] obiceretur V S n71
114 sunt] sint S n72

# Interoperability: the issue

- Processing
  - Search
  - Collation
  - **NLP (lemma, PoS etc.)**
  - Statistics (dist. reading)

uenenú

- uenenū

- uenenum

# Interoperability: the issue

- Processing
  - Search
  - Collation
  - NLP (lemma, PoS etc.)
  - **Statistics (dist. reading)**

uenenū

uenenum

# Interoperability: the issue

- ***My focus: European Medieval handwriting***
  - ...and early print (imitating handwriting)

# Unicode (TEI's recommendation)

- Solution for new digital texts

- Not enough for pre-modern writing systems
  - Allographs
    - ſ (U+017F)  /  s (U+0073; ASCII 115)
    - Have I told the computer that they correspond to each other (variants of grapheme <s>)?
  - Ligatures
    - & (U+0026; ASCII 38)
    - Have I encoded that it is equivalent to "e + t" in that MS?
  - Grapheme set
    - u (U+0075; ASCII 117)
    - Have I encoded whether it "covers" (or not) <u> *and* <v>?

# Manual normalized transcription

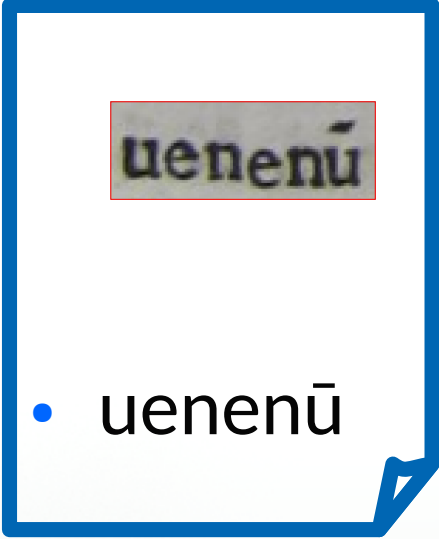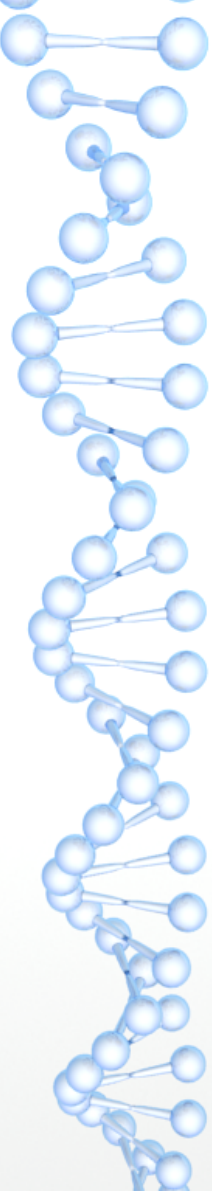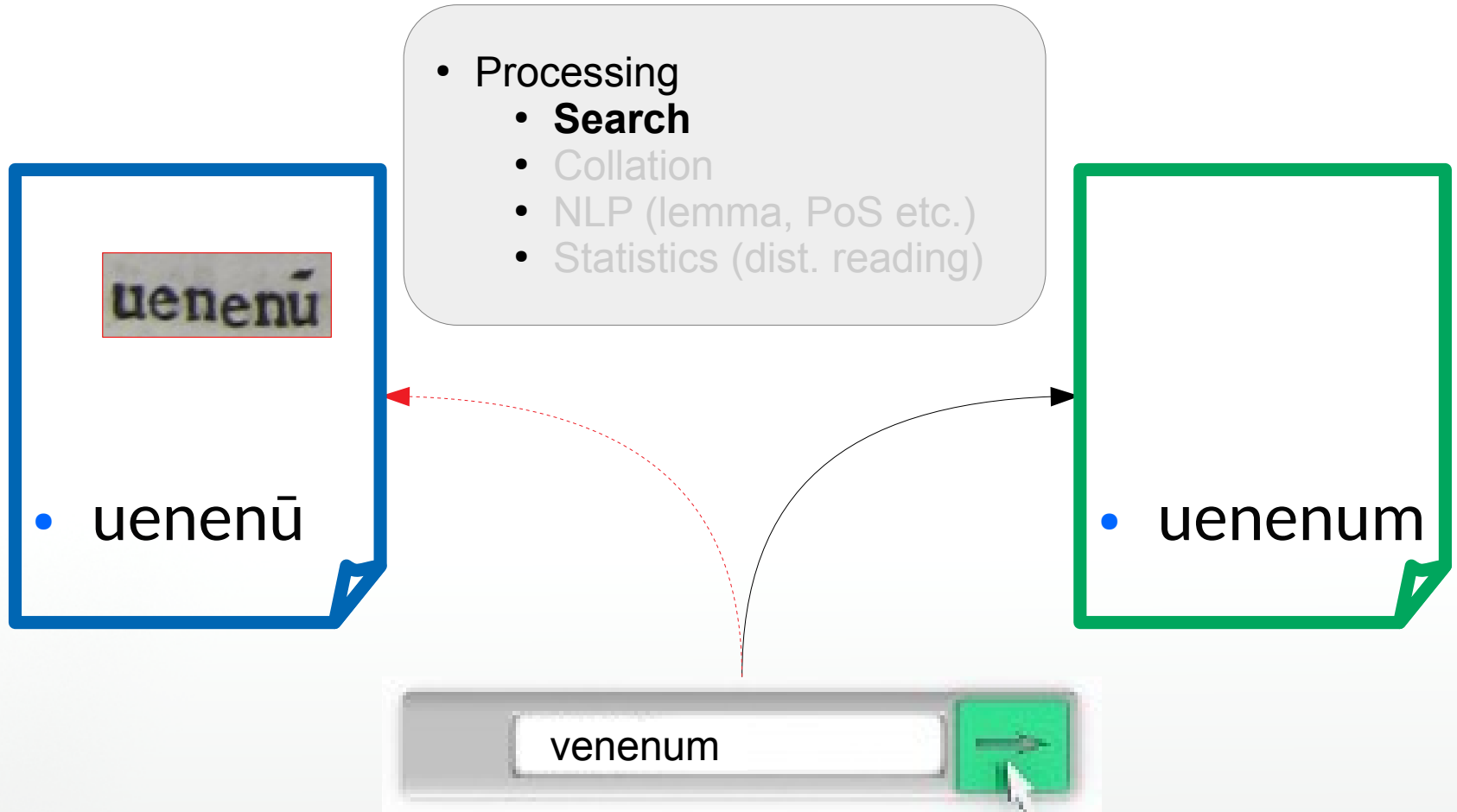uenenú

- uenenū

Diplomatic

- Visualization
- ~~Processing~~

# Manual normalized transcription

- venenum



- uenenū

**Normalized**

- Processing
  - Search
  - Indexing
  - Collation
  - NLP (lemma, PoS etc.)
  - Statistics (distant reading)...

**Diplomatic**

- Visualization
- ~~Processing~~

# Manual normalized transcription

- venenum

Not generated by computer: repeated manual transcription

uenenú

- uenenū

**Normalized**
- Processing
  - Search
  - Indexing
  - Collation
  - NLP (lemma, PoS etc.)
  - Statistics (distant reading)...

**Diplomatic**
- Visualization
- ~~Processing~~

# Project-specific pre-processing

- Collation: pre-processing to normalize the text
  - Or instruct collation software to ignore *some* discrepancies
- Same with search/indexing/NLP/statistics software
- Disposable home-made solutions
  - ...for a shared issue

# Interoperability through modelling

WANTED

DEAD OR ALIVE
REWARD
$4,000

- Documenting project-specific modelling (transcription) and normalization practices
  - What am I transcribing (graphemes/allographs)? How?
  - What corresponds to what (f/s)?
- Documentation
  - In English prose (specifications for other programmers)
  - Formal (software code, tables)

# Interoperability through modelling

WANTED

DEAD OR ALIVE
REWARD
$4,000

- Documenting project-specific modelling (transcription) and normalization practices
  - What am I transcribing (graphemes/allographs)? How?
  - What corresponds to what (f/s)?
- Documentation
  - In English prose (specifications for other programmers)
  - Formal (software code, tables)

Example ahead!

# Orlandi's table of signs

OCR/HTR (witness A)

Manual (selective) transcription (witness B)

# Orlandi's table of signs



**Allographic transcription**

Vnder τhis Casτle

under this caſtle

OCR/HTR (witness A)

Manual (selective) transcription (witness B)

# Orlandi's table of signs



**Allographic transcription**

| | |
|---|---|
| Vnder ᴛhis Casᴛle | under this caſtle |

OCR/HTR (witness A)

Manual (selective) transcription (witness B)

# Orlandi's table of signs



Unicode characters

Allographic transcription

Vnder ⊤his Cas⊤le

under this ca∫tle

OCR/HTR
(witness A)

Manual (selective)
transcription
(witness B)

# Orlandi's table of signs



**Allographic transcription**

Vnder тhis Casтle

under this caſtle

OCR/HTR (witness A)

Manual (selective) transcription (witness B)

# Orlandi's table of signs



Allographic transcription

Vnder ꞇhis Casꞇle

under ꞇhis caſtle

OCR/HTR (witness A)

Manual (selective) transcription (witness B)

# Orlandi's table of signs



**Allographic transcription**

Vnder τhis Casτle

under this caſtle

OCR/HTR (witness A)

Manual (selective) transcription (witness B)

# Orlandi's table of signs



Allographic transcription

Vnder ꞇhis Casꞇle

under this caſtle

OCR/HTR (witness A)

Manual (selective) transcription (witness B)

# Orlandi's table of signs



| Gr | Allogr |
|----|--------|
| s: | s |
| t: | τ \| ℇ \| √ |
| u: | u \| V |

| Gr | Allogr |
|----|--------|
| s: | s \| ſ |
| t: | t |
| u: | u |

**Allographic transcription**

Vnder τhis Casτle

under this caſtle

OCR/HTR (witness A)

Manual (selective) transcription (witness B)

# Orlandi's table of signs



Graphematic layer

unter dem schloss

unter dem schloss

```
Gr    Allogr
s:    s
t:    τ | ε | √
u:    u | ᴠ
```

```
Gr    Allogr
s:    s | ſ
t:    t
u:    u
```

- **Not** hard-coded in the project-specific normalization practice/software
- **But** documented and formalized

Allographic transcription

Vnder ᴛhis Casᴛle

under thiſ caſtle

OCR/HTR (witness A)

Manual (selective) transcription (witness B)

# Orlandi's table of signs

**Graphematic layer**

unter dem schloss

unter dem schloss

**(More) interoperability**

- Processing
  - Search
  - Indexing
  - Collation
  - NLP (lemma, PoS etc.)
  - Statistics (dist. reading)

```
Gr      Allogr
s:      s
t:      τ | ε | √
u:      u | V
```

```
Gr      Allogr
s:      s | ſ
t:      t
u:      u
```

**Allographic transcription**

Vnder τhis Casτle

under this caſtle

- Visualization
- Processing

OCR/HTR (witness A)

Manual (selective) transcription (witness B)

# Interoperability through modelling

WANTED  DEAD OR ALIVE  REWARD $4,000

- Documenting project-specific modelling (transcription) and normalization practices
  - What am I transcribing (graphemes/allographs)? How?
  - What corresponds to what (f/s)?
- Documentation
  - In English prose (specifications for other programmers)
  - Formal (software code, tables)

# Interoperability through modelling

WANTED    DEAD OR ALIVE
REWARD
$4,000

- Scholarly discussion on **modelling**
  - What is a grapheme? What is an allograph?
- **Shared** models
  - Grassroot approach, from discussion
- Reusable **software** libraries

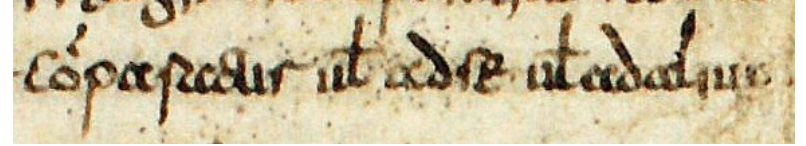# Graphemes/allographs: the commutation test

System

Comparatur vel ad se vel ad alium
*He is compared to himself or to another*



\<s\>
\<t\>

Text

- cȯparaƐur   uł adſe   uładalium

\<x\>
\<y\>
\<z\>

# Graphemes/allographs: the commutation test

# Graphemes/allographs: the commutation test

«τ»

&lt;s&gt;
&lt;t&gt;

• c̓paraƐur      uɫ adʄe   uɫadalium

«√»

&lt;x&gt;
&lt;y&gt;
&lt;z&gt;

**Substitution**:
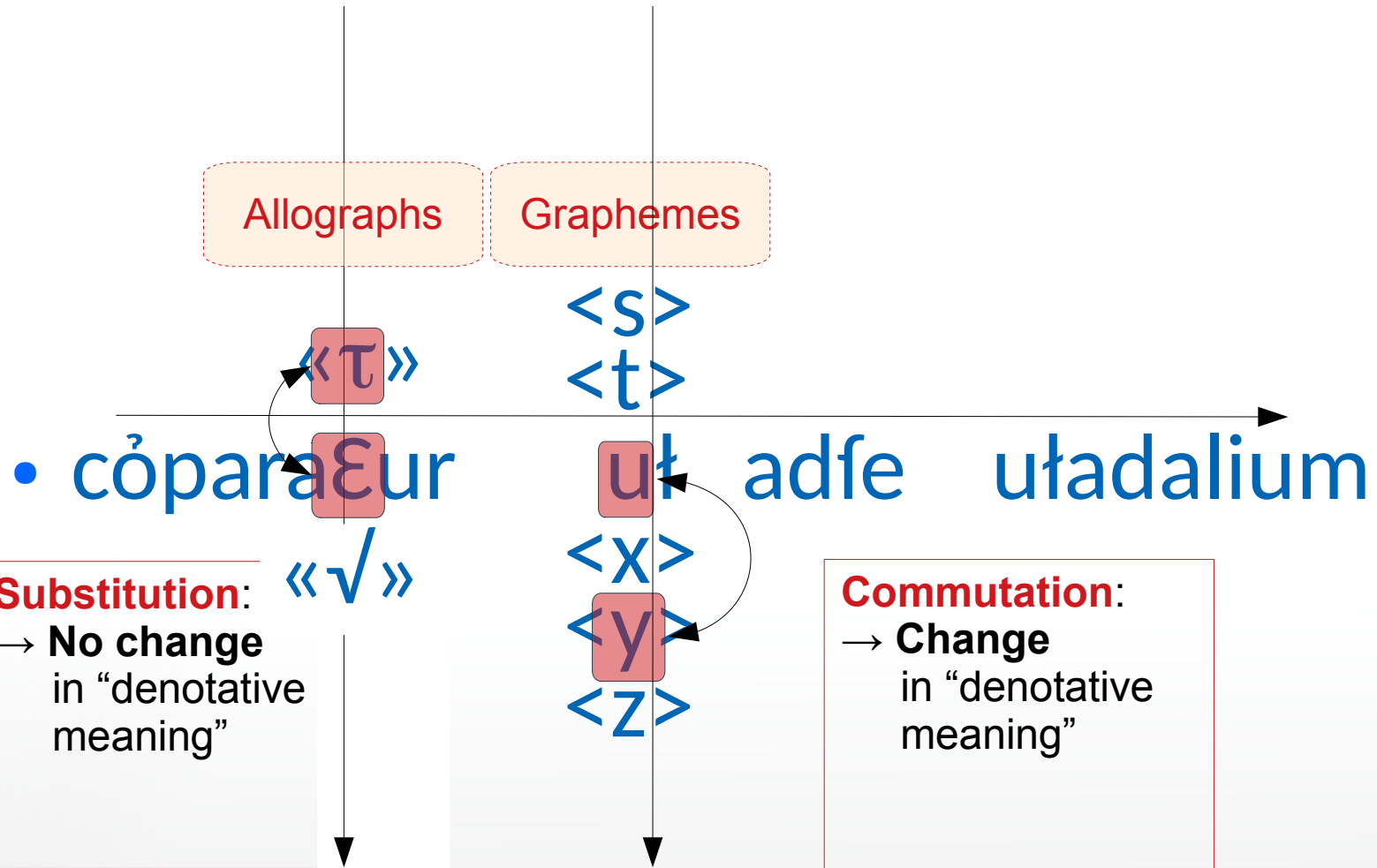→ **No change**
in "denotative meaning"

**Commutation**:
→ **Change**
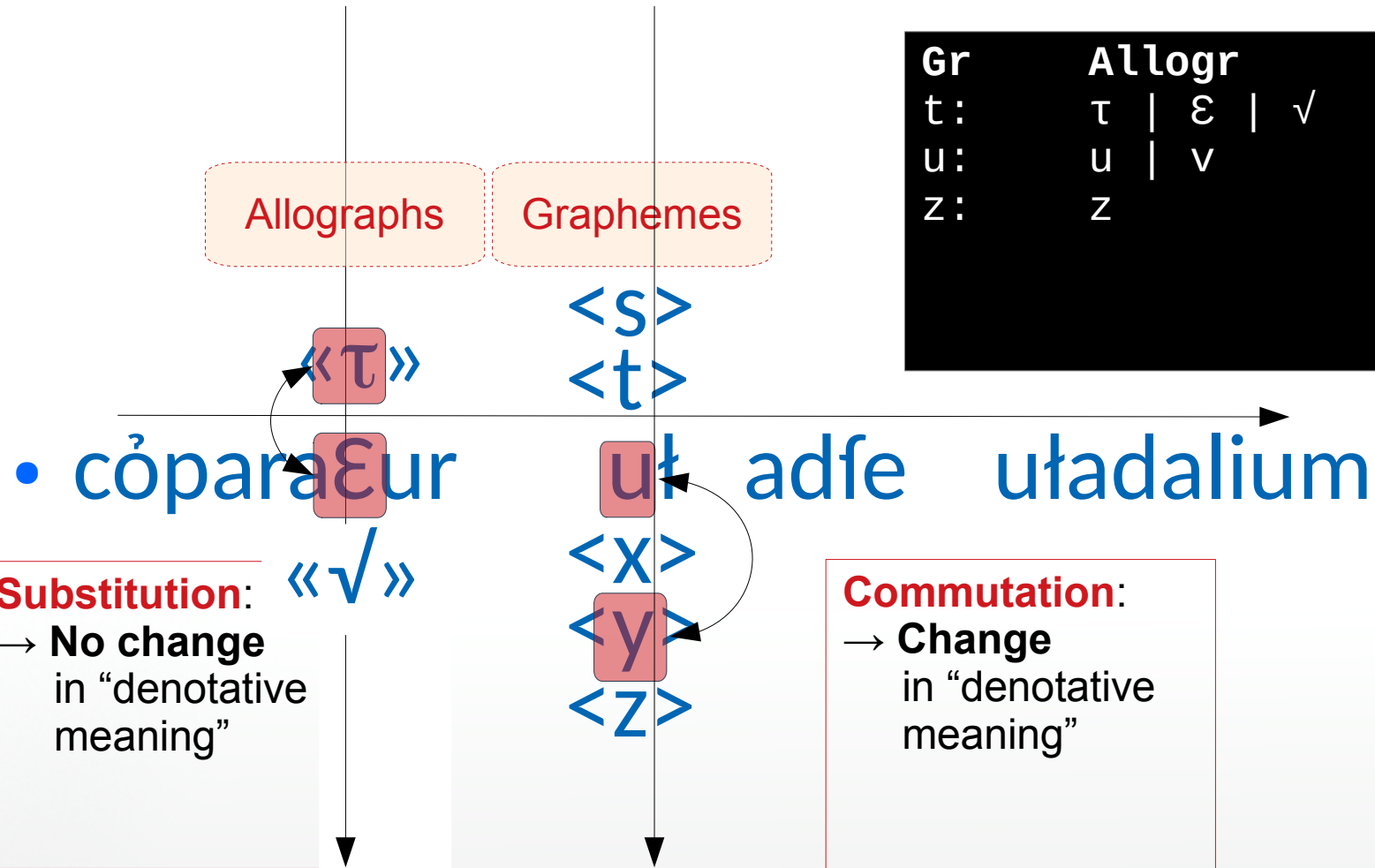in "denotative meaning"

# Graphemes/allographs: the commutation test

# Graphemes/allographs: the commutation test

| Gr | Allogr |
|----|--------|
| t: | τ \| Ɛ \| √ |
| u: | u \| v |
| z: | z |

Allographs

Graphemes

‹s›
‹t›
«τ»
ċoparaƐur    uł adſe   uładalium
«√»

‹x›
‹y›
‹z›

**Substitution**:
→ **No change**
in "denotative meaning"

**Commutation**:
→ **Change**
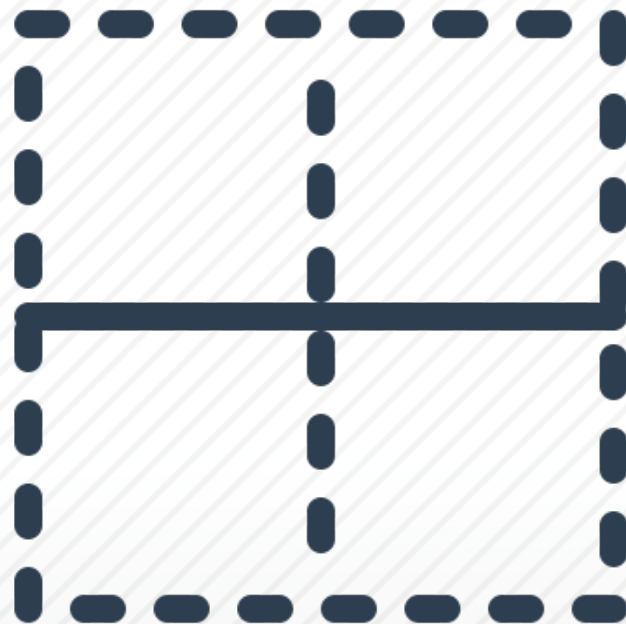in "denotative meaning"

# Graphemes / allographs: what to transcribe?

- Whatever our project needs!
  - Based on its scientific interests
  - (…and on time / money)
- But: *declaring* what we are transcribing
  - Based on a formal distinction, e.g., between graphemes / allographs
  - Documenting it
    - English prose
    - Formal (software, tables…)

Good practices

# Good practices

- See section Seminar readings on GitHub
    - Human-readable tables of signs
    - Machine-readable tables of signs

# Exercise

- See section Exercise on GitHub

# Outline

- **Interoperability**
  of digital scholarly editions (DSEs)
  based on diplomatic transcriptions

  - The issue

  - Current solutions

  - Interoperability through modelling

- Orlandi's **table of sign**

- **Graphemes/allographs**