

Including XML Markup in the Automated Collation of Literary Text

Elli Bleeker, Bram Buitendijk, Ronald Haentjens Dekker, Astrid Kulsdom

*R&D - Huygens Institute for the History of the Netherlands
KNAW Humanities Cluster*



@ellibleeker
@ronald_dekker
@bram_buitendijk

Context: XML in the Humanities

Text modelling outside OHCO model: TEI (Text Encoding Initiative)

Textual phenomena (variation, overlap, etc.)

Focus: computational philology and textual variation

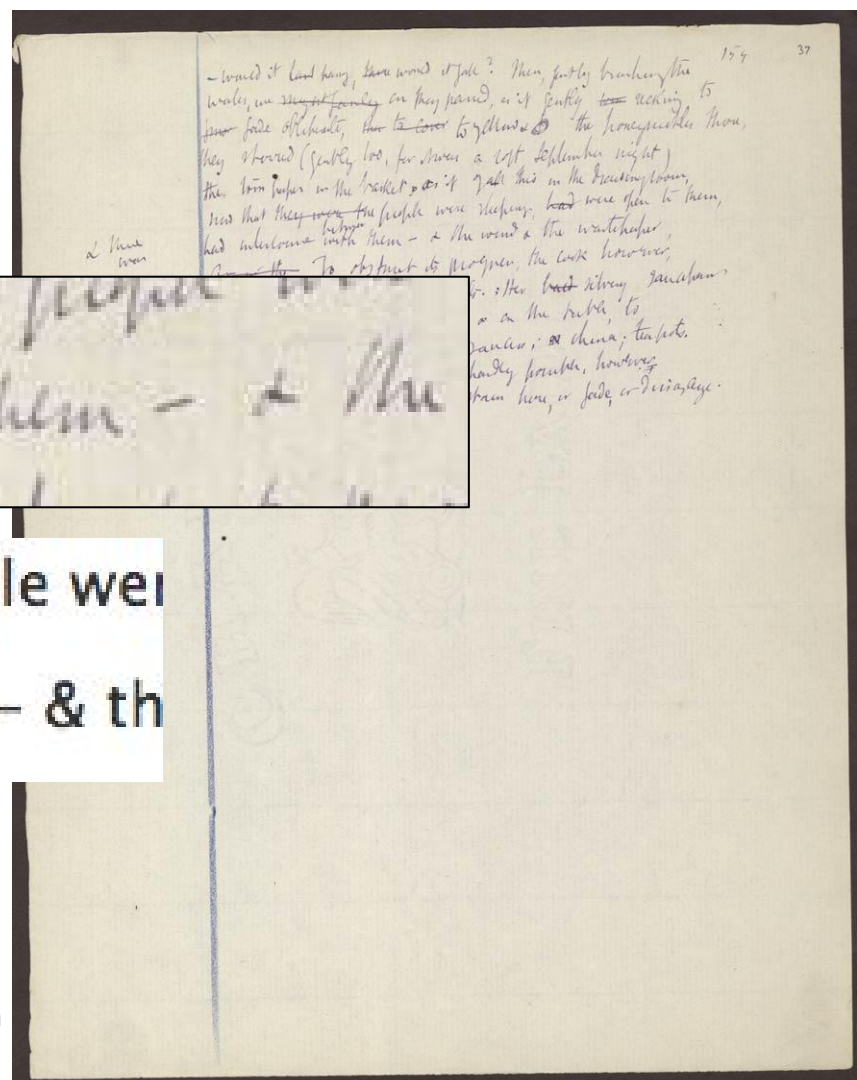
"Collation": the comparison of two or more text versions

Problem Statement

How can we make optimal use of the potential of XML for computational philology?

1. Modelling textual variation in TEI
2. Hypergraphs and hypergraph merging
3. Visualise/export results

Open Variants

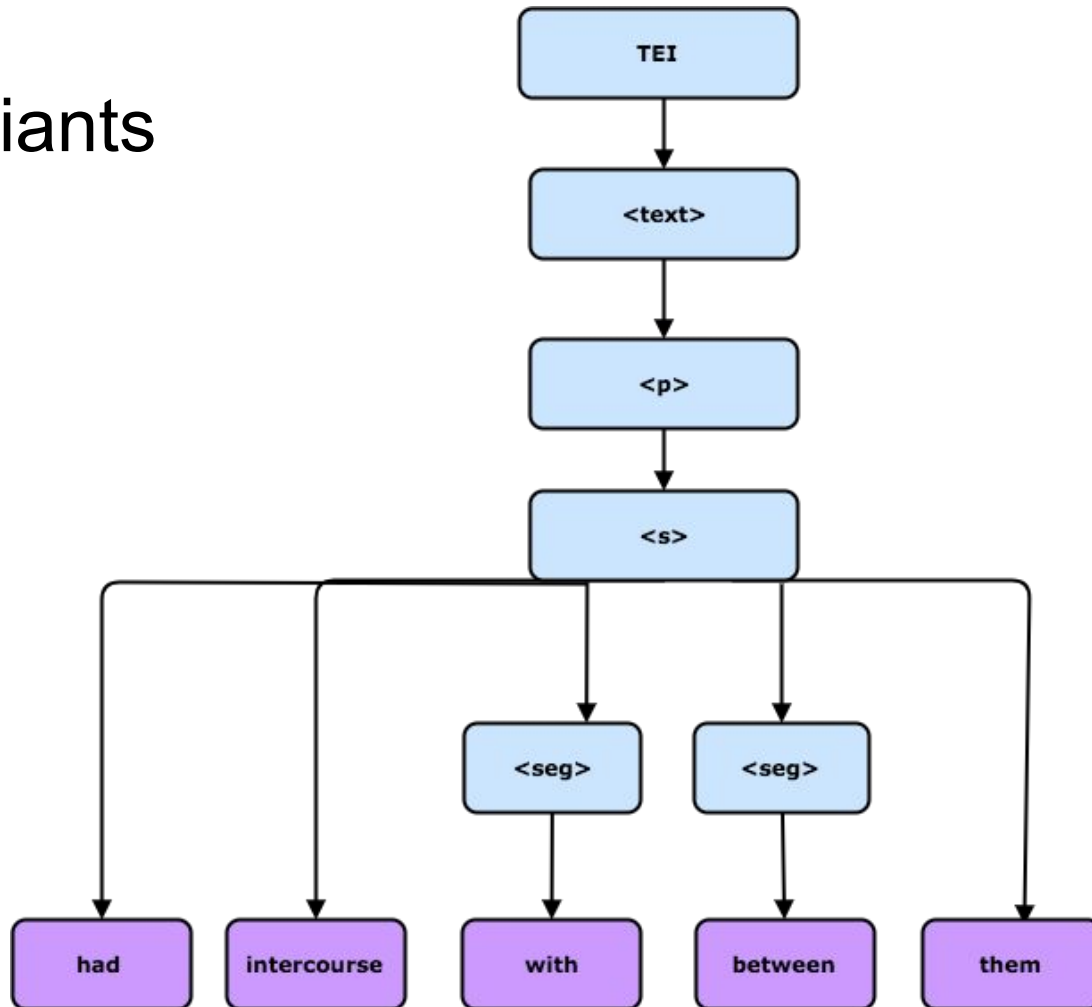


now that they were the people we
had intercourse with them -- & th

Open Variants

```
<?xml version="1.0" encoding="UTF-8"?>
<TEI>
  <text>
    <p>
      <s>were open to them, had intercourse <seg type="alternative" xml:id="alt1">with</seg>
        <seg type="alternative" xml:id="alt2">between</seg> them</s>
    </p>
  </text>
</TEI>
```

Open Variants



Simultaneity

*folding the house
in the land
a board creaks.*

**[] pulse. & then gently mak muffling & folding the house again.
on the landg**
Only suddenly & for no perceptible reason, a board creaks; &

167⁵
this air of simple integrity, for whose sake is no strife or contention
no compromise, it seems as if truth were there, undraped,
itself, undisturbed.

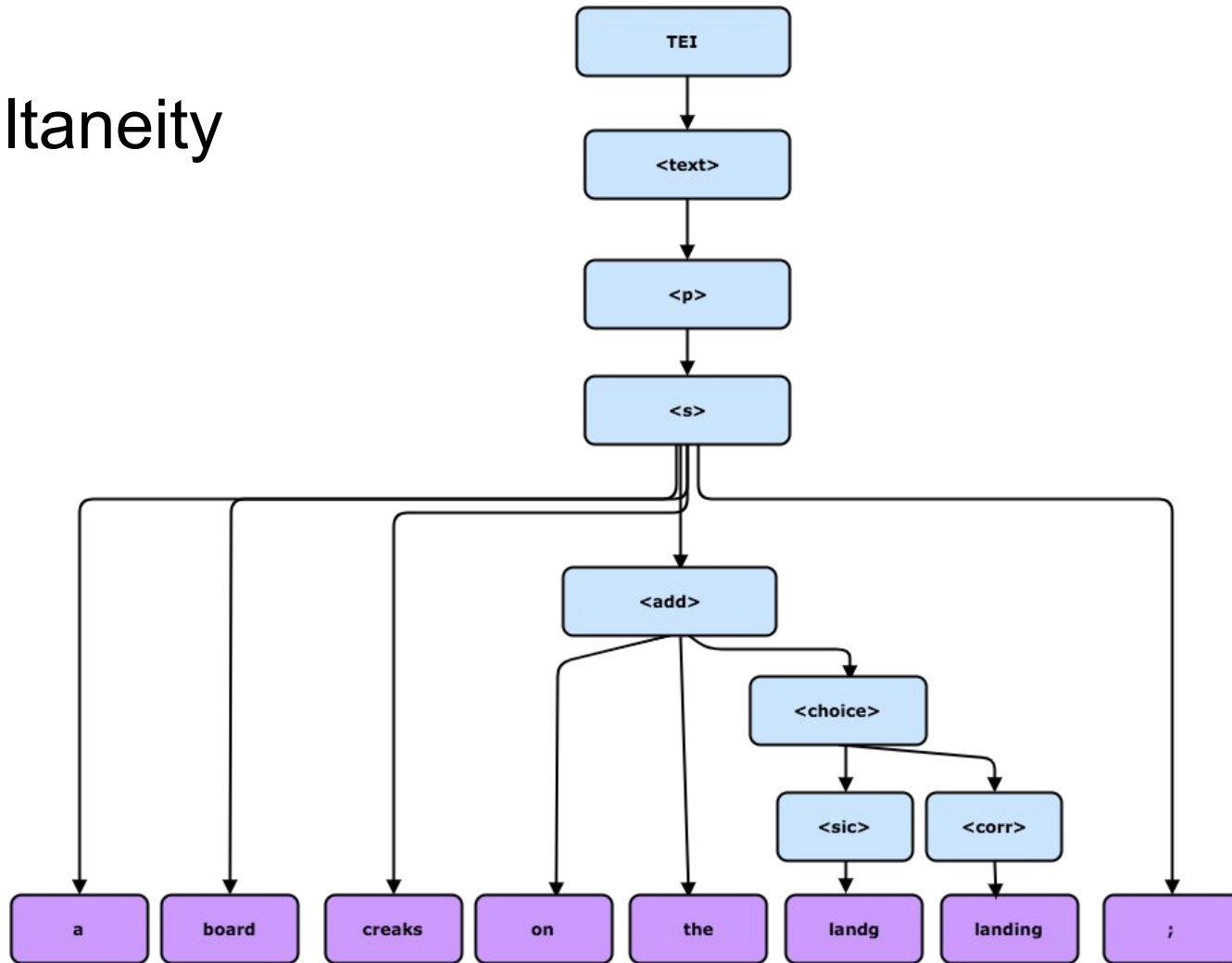
So ~~the~~ what its nothing it seems can break that fear
image: or comfort that innocence; or disturb the lonely way
the splendence which day after day work a week &
glee might lie upon the face, whose gaze for the
lusts & the damp sea air & weaves with
tumbling cry of rock on the ~~fallen~~ ~~fallen~~ ~~fallen~~ ~~fallen~~ ~~fallen~~ ~~fallen~~ ~~fallen~~
coming from the sea, or
for awhile in the field of some some have rising
in the back, whether some there: some
& they gently mak muffling - folding the house again.
ly, & for no perceptible reason, a board creaks:

in against the weather; the brave looses itself as
of centuries in one second a giant avalanche detach-
entirely & where the weather are fold the
awakened, & whether by it as it came
edwin's few & in early it speaks; had finally bitten
through some ice & loosened the shafts hold when the
shell. at the same moment too, - almost more
purposefully, with a noise like the grating of iron
& the

Simultaneity

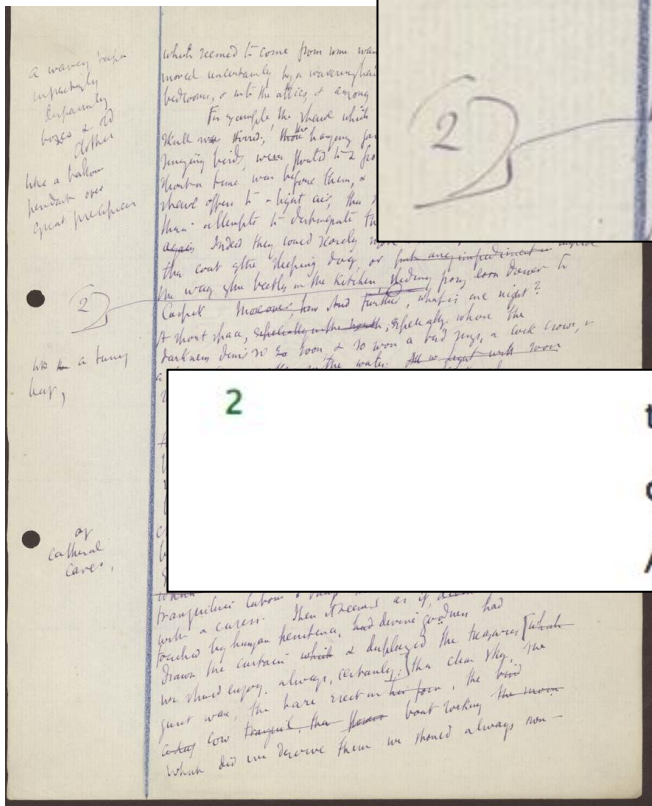
```
<?xml version="1.0" encoding="UTF-8"?>
<TEI>
  <text>
    <p>
      <s>a board creaks <add>on the <choice><sic>landg</sic><corr>landing</corr></choice></add>;</s>
    </p>
  </text>
</TEI>
```


Simultaneity



Witness 1

Structure



the coat of the sleeping dog, or ~~put any impediment~~
the way of the beetles in the kitchen sliding from ~~corn~~ dresser to
Carpet. Moreover, how And Further, what is one night?
A short space, especially in the north; especially where the

the way of the beetles in the kitchen, sliding from ~~corn~~ dresser to
carpet. Moreover, how And Further, what is one night?
A short space, especially in the north; especially where the

Witness 2

Structure

earth seemed ruining and washing away in water.

III

But what after all, is one night? A short space,
especially when the darkness dims so soon, and so soon a bird

earth seemed ruining and washing away in water.

III

But what after all, is one night? A short space,
especially when the darkness dims so soon, and so soon a bird

(6)

they spread their garments, they rose up and
the waves rose and through the house there
sullen wave of doom which curled and crashed
earth seemed ruining and washing away in water.

III

But what after all, is one night? A short space,
especially when the darkness dims so soon, and so soon a bird
sings, a cock crows, or a faint green quick
leaf, in the hollow of the wave. Night, however
night. The winter holds a pack of them in its
equally, evenly, with indefatigable fingers.

they darken. Some of them hold aloft clear
brightness. The autumn trees, ravaged as the
flash of tattered flags kindling in the gloom
caves where gold letters and marble pages descend
battle and how bones far away bleach and burn
The autumn trees gleam in the yellow moonlight
harvest moons, the light which mellow the one
smooths the stubble, and brings the wave lapping

It seemed now as if, touched by human pen
toil, divine goodness had drawn the curtain
it, single, distinct, the hare erect, the wave
boat rocking, which, did we deserve them should be ours always.

But alas - divine goodness, twitching the cord, draws the curtain:
it does not please him; he covers his treasures in a drench of

Structure

Witness 1

```
<?xml version="1.0" encoding="UTF-8"?>
<TEI>
  <div type="chapter" n="1">
    <p>
      <s>dresser to carpet.</s>
      <add place="margin"><head>2</head></add>
      <s>Further, what is one night?</s>
    </p>
  </div>
</TEI>
```

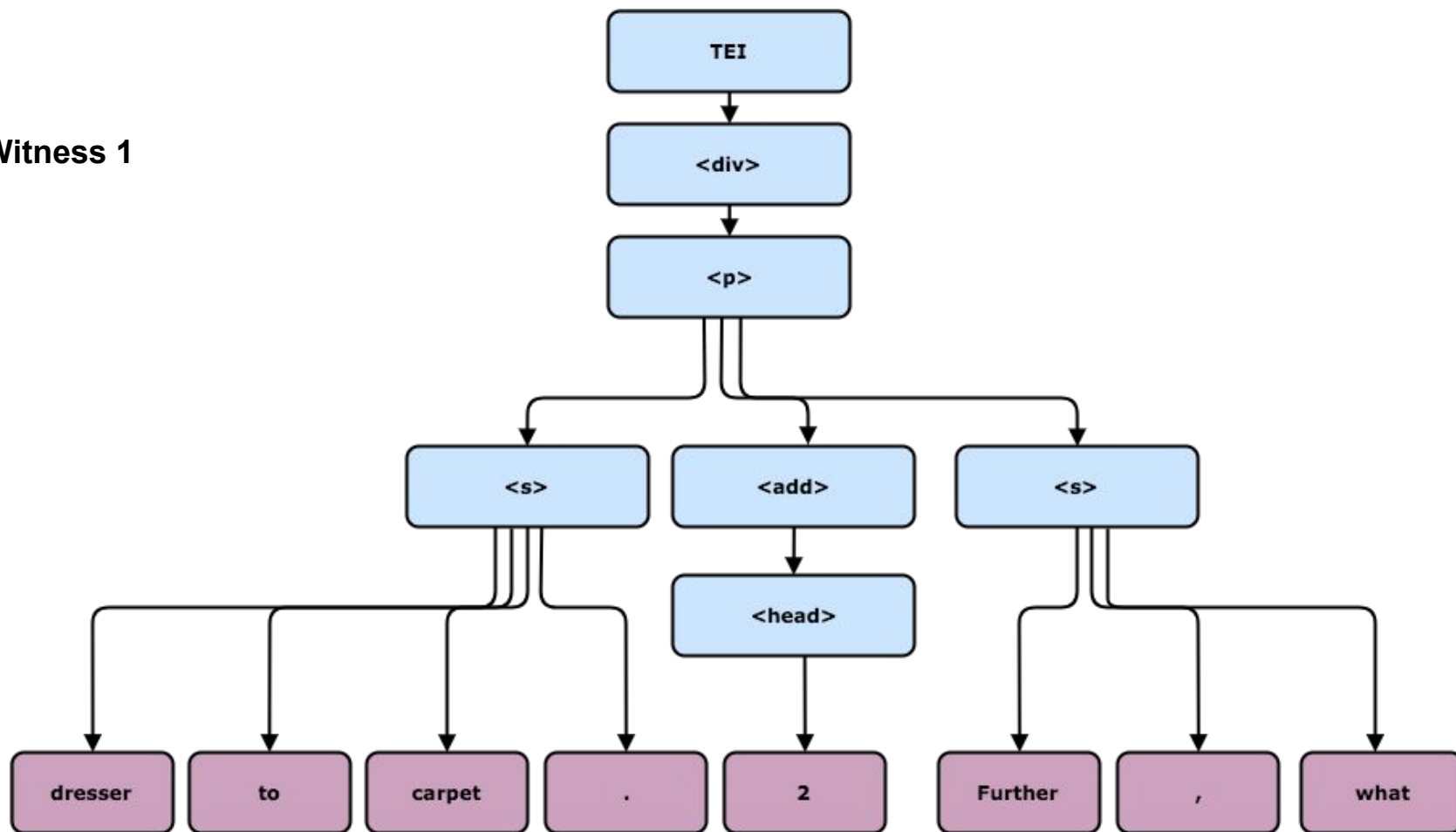
(simplified) TEI-XML transcription of manuscript p.157

Witness 2

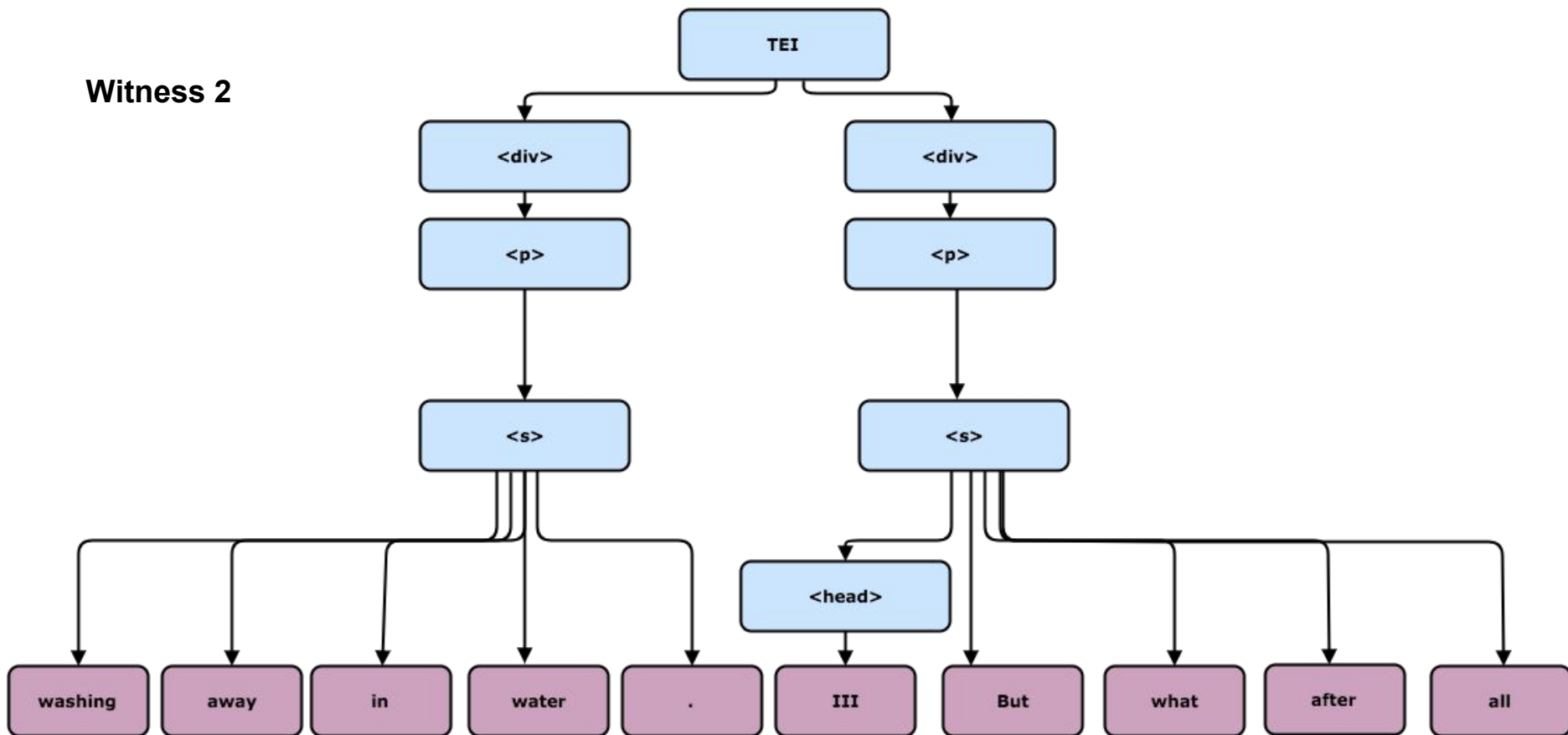
```
<?xml version="1.0" encoding="UTF-8"?>
<TEI>
  <div type="chapter" n="2">
    <p><s>washing away in water.</s></p>
  </div>
  <div type="chapter" n="3">
    <head>III</head>
    <p><s>But what after all, is one night?</s></p>
  </div>
</TEI>
```

(simplified) TEI-XML transcription of typescript p.5

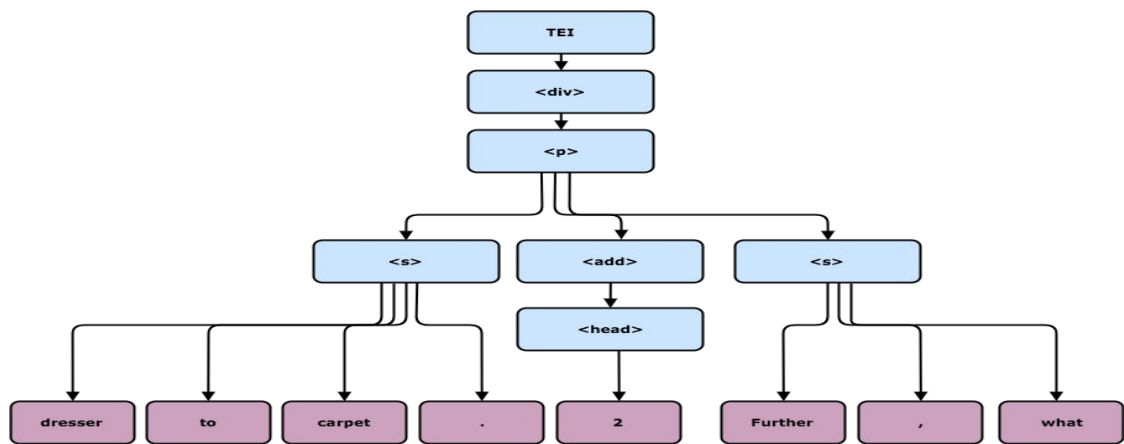
Witness 1



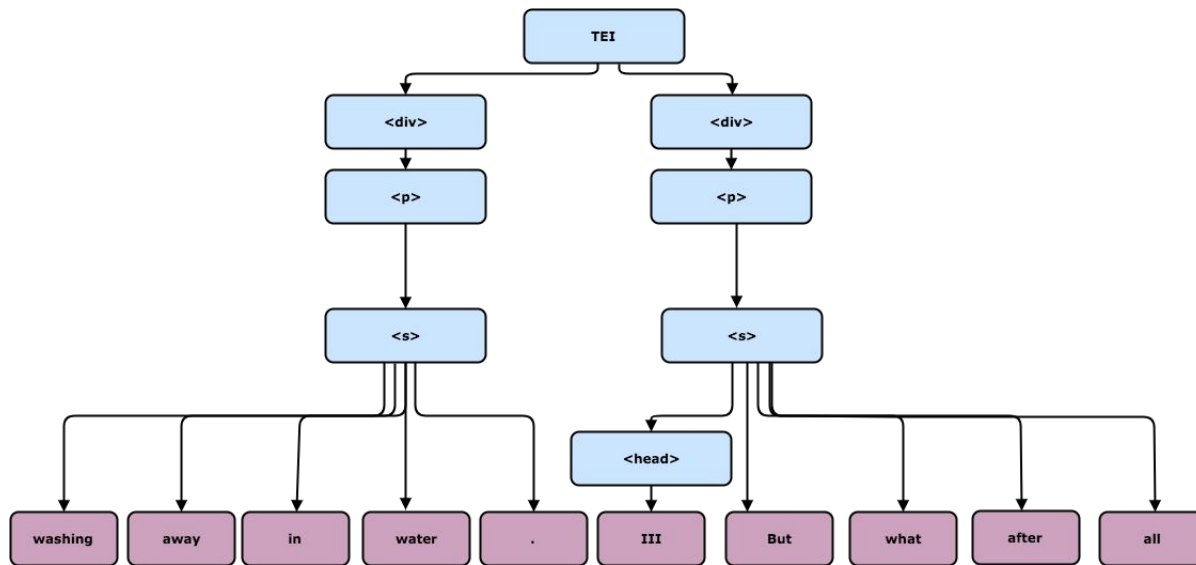
Witness 2



Witness 1



Witness 2



Modelling Texts in the Humanities

Beyond OHCO:

- Intradocumentary variation → multiple paths through one text, encoded with TEI tags
- Structure → comparing different documents with different structures results in conflicting hierarchies

Conditions

Different types of XML data:

- fully ordered
- unordered

Text-centric XML of literary texts (in short: "TEI text") = partially ordered data

Functional Requirements

Processing and analysing TEI text in a native way

Automated collation: finding the minimum set of changes needed to turn one document into the other.

Outset:

- documents are partially ordered
- treat structure and text equally
- compare more than 2 documents

Technical Requirements

Data model:

- schema-aware
- store multiple hierarchies

Processing:

- find the smallest amount of differences between two TEI texts
- distinguish between words, punctuation, markup
- respect non-linearity of TEI text

HyperCollate

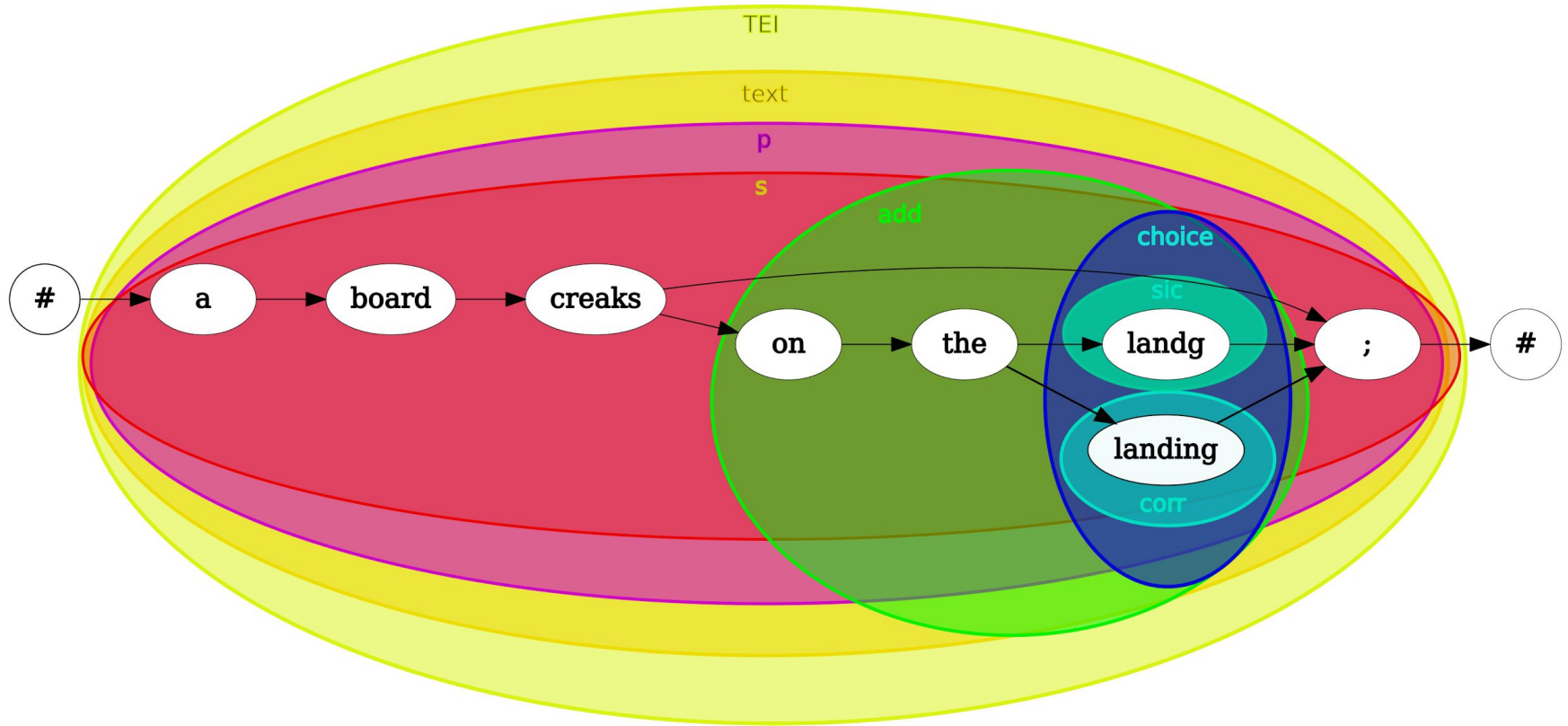
HyperCollate = automated collation tool that uses a hypergraph model for variation

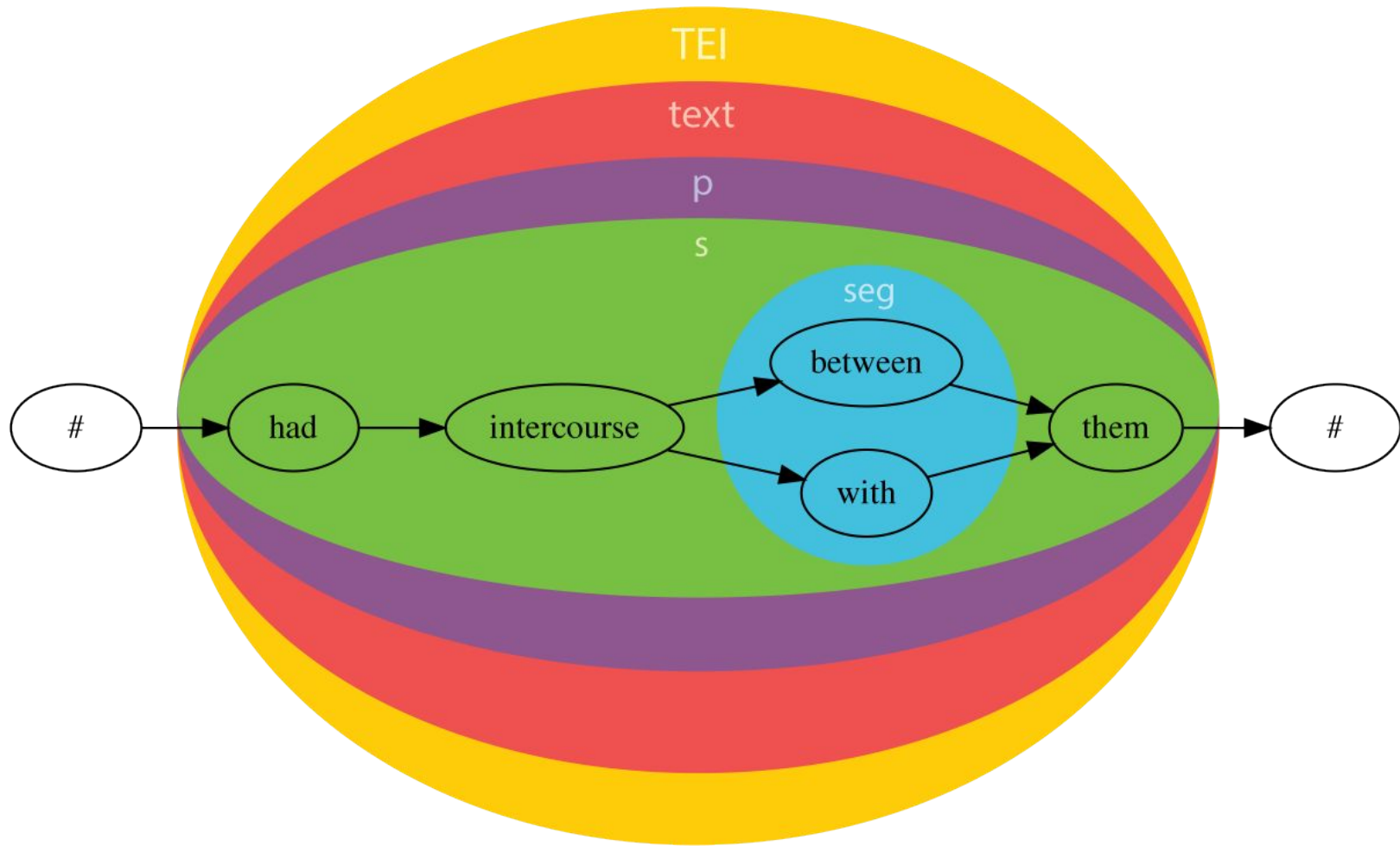
Hypergraph: nodes, hyperedges (one-to-one; one-to-many; many-to-many)

HyperCollate processes complex features of text, including:

- overlapping hierarchies
- intradocumentary variation

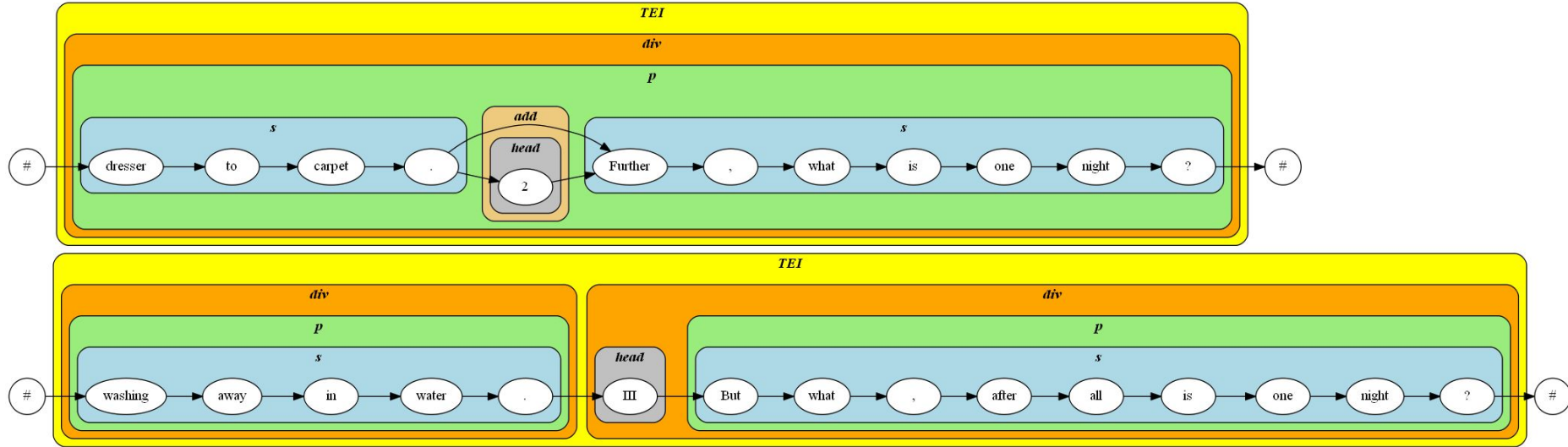
Hypergraph Model for Textual Variation



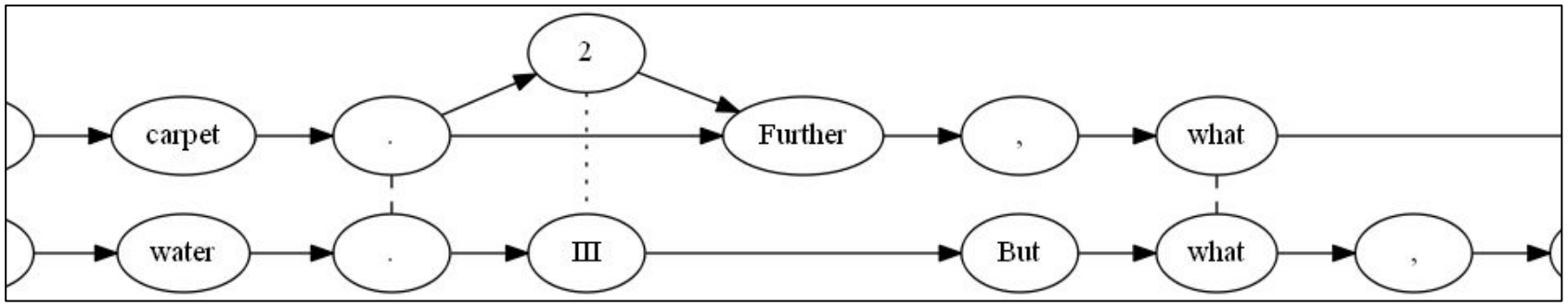
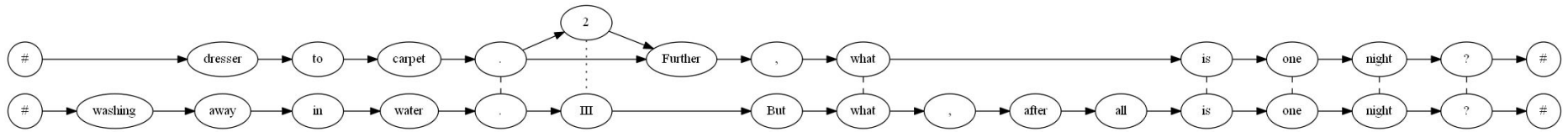


Steps of HyperCollate

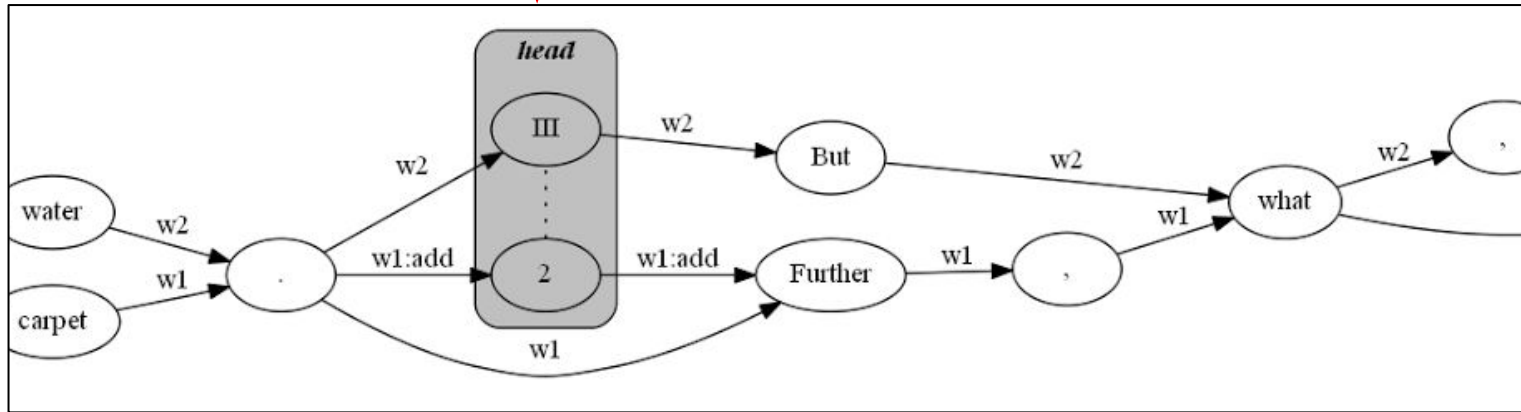
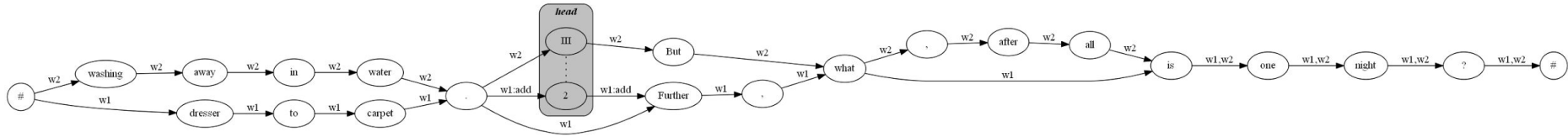
1. Convert XML tree of each witness into a variant hypergraph
2. Align two variant hypergraphs
3. Merge two hypergraphs in one collation hypergraph
4. Repeat in case of >2 witnesses
5. Visualise/export collation hypergraph



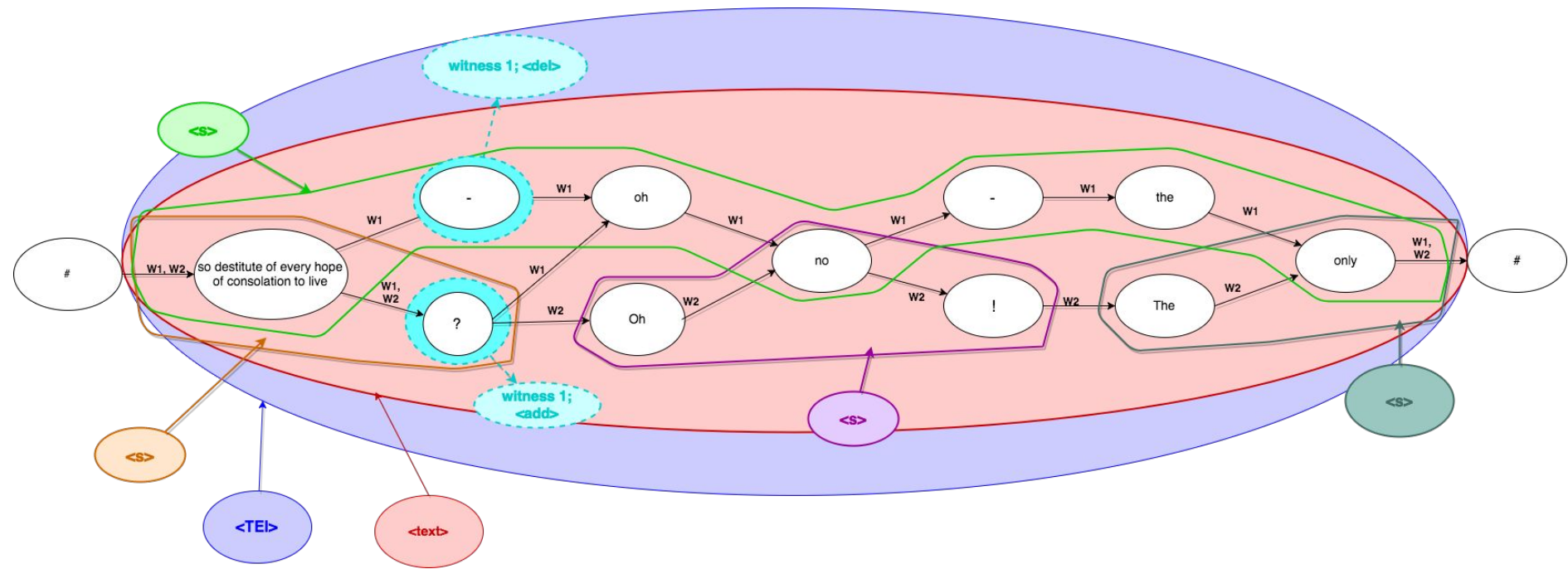
1. Transform TEI text witnesses into variant hypergraphs



2. Align the variant hypergraphs



3. Merge variant hypergraphs into one collation hypergraph



Conclusion

TEI text, being partially ordered data, is especially challenging for processing.

Requirements for analysis of textual variation:

- process multiple paths (i.e. be schema-aware)
- store multiple hierarchies

A hypergraph data model allows us to make optimal use of XML's potential for humanities research.

Elli Bleeker, Bram Buitendijk, Ronald Haentjens Dekker, Astrid Kulsdom

R&D - Huygens Institute for the History of the Netherlands

KNAW Humanities Cluster



@ellibleeker

@ronald_dekker

@bram_buitendijk