

# Lyric Alignment on Plainchant Manuscripts

Tim de Reuse

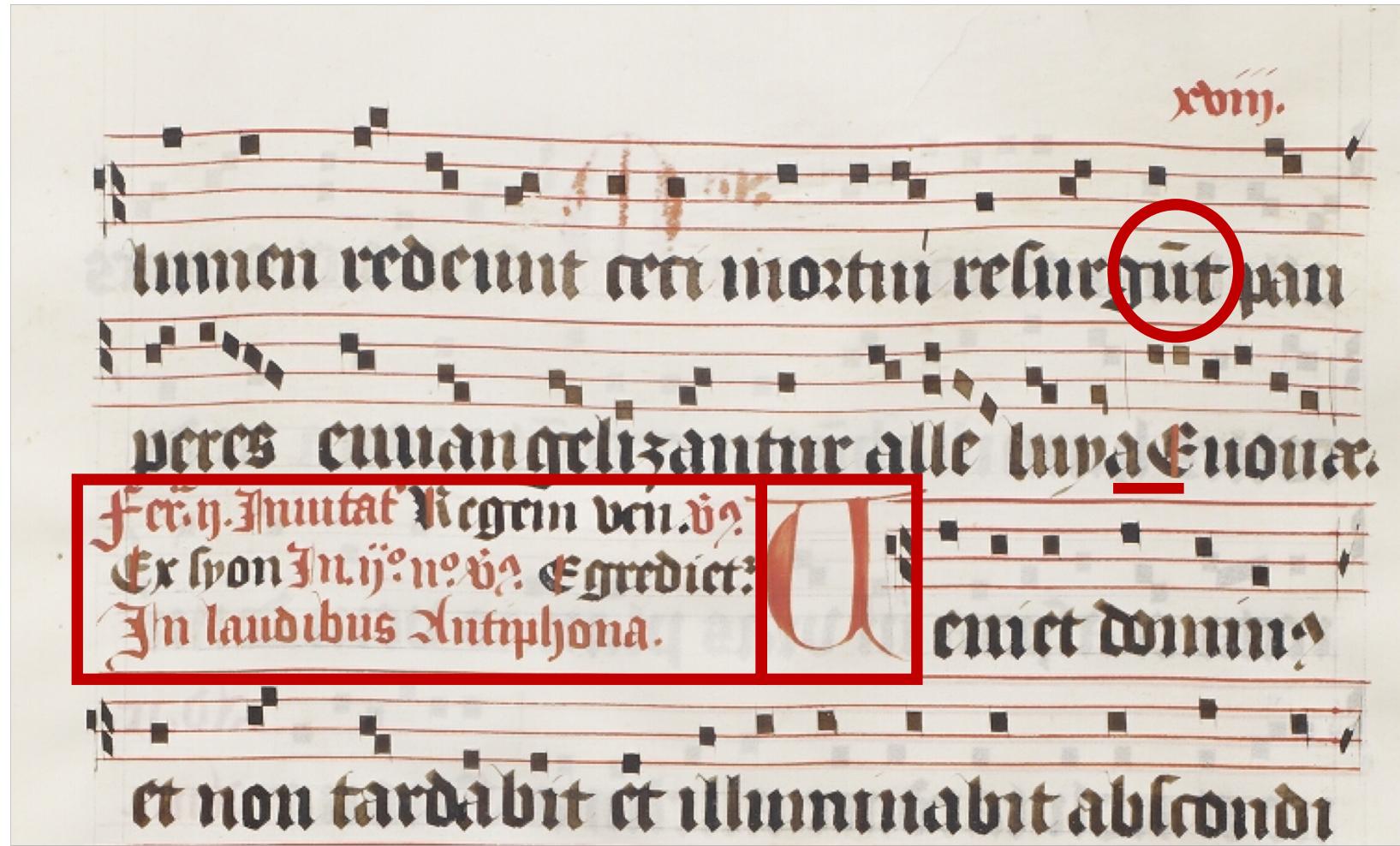
SIMSSA Workshop XVII

1/12/2018

# Associating Neumes to Syllables of Text

- Must locate text first
- Optical Character Recognition (OCR) on handwritten sources is difficult!
- Existing transcripts do not indicate *where* text is on page





Non-musical Text

Abbreviations

Inconsistent  
Spacing

Ornamental Letters

**TRANSCRIPT:** lumen redeunt ceci mortui resurgunt pauperes  
euvangelizantur alleluia Euouae Veniet dominus et non tardabit et  
illuminabit abscondi (...)

# Aligning to Imperfect OCR

- Run OCR on handwritten source
  - Find out where *most* of the characters are
  - Expect lots of errors
- Compare OCR result with correct transcript
  - Correct errors in the OCR and find all characters
  - Identify non-musical text, abbreviations, etc.

# Training the OCR

- ~650 lines of text transcribed for training a model
  - 42 Pages of the Salzinnes Antiphonal
  - That's not a lot for handwritten OCR!
  - Using OCropus open-source OCR system
- ~80% per-character error rate after 24 hours training
  - Good enough!

## **TRANSCRIPT:**

lumen redeunt ceci mortui resurgunt pauperes euvangelizantur alleluya Euouae

## **OCR:**

lmmen redeunt ceci mortiiui resurgūt pan yeres eyugelisantur alle luya Euonae

# Global Sequence Alignment

- Given two sequences, edit them to make them “match up”
  - Cost minimization: as few gaps as possible for as many matches as possible

  Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor  
  incididunt ut labore et dolore magna aliqua.

+

  Lorem ipsum dollllllor acsit amet, consectetur di.s elit,, eiusmmd tempods  
  incididunt ut lb ore etmagna aliqua.

↓

  Lorem ipsum do\_\_\_\_\_lor \_\_\_\_sit amet, consectetur adipiscing elit, sed do  
  Lorem ipsum dollllllor acsit amet, consec\_\_\_\_tur \_di\_.s\_\_\_\_ elit,\_\_\_\_\_,

  eiusmod tempo\_\_\_\_r incididunt ut labore et dolore magna aliqua  
  eiusmmd tempods\_\_\_\_r incididunt ut lb ore et\_\_\_\_\_ magna aliqua

# Aligning the Transcript to OCR Output

## TRANSCRIPT:

lumen redeunt ceci mortui resurgunt pauperes euvangelizantur alleluya Euouae

## OCR:

lmmen redeunt ceci mortiiui resurgūt pan yeres eyugelisantur alle luya Euonae

## ALIGNMENT:

lumen redeunt ceci mort\_\_ui resurgunt pauperes

lmmen redeunt ceci mortiiui resurgu\_t panyeres

euvangelizantur alle\_luya Euouae

e\_\_yugelisantur alle luya Euonae.

The OCR gets characters wrong, but the incorrect character  
is often aligned to the correct one in the transcript

lumen redeunt ceci mortui resurgunt pa

peres euangelizantur alle luya Euouae.

Fecit. Invitat Regem vestrum.  
Ex syon inviso nova Egressus.  
In laudibus Antiphona.

**U**eni et dominus

et non tardabit et illuminabit abscondi

lumen redeunt ceci mortui resurgunt pauperes euangelizantur alle luya Euouae  
lumen redeunt ceci mortui resurgunt pauperes euangelizantur alle luya Euouae. cx.n nni

\_\_\_\_\_  
Ven\_\_\_\_\_ iet dominus et non tar  
tat segeim ven.vus. In landibus Antiphona In laudibus Antiphona emet dominus et non tar

A long gap: The OCR found non-musical text not present in the transcript,  
and the sequence alignment identified it as such

I u m e n r e d e u n t cec̄t̄ m o r t u i r e s u r g u n t p a u  
**lumen redemit arc̄i mortuū resurgūt̄ pa**

p e r e s e u v a n g e l i z a n t u r a l l e l u y a E u o u a e  
**p̄c̄es euangelizantur alle luva Euouat̄**  
Fer̄ij Inuitat̄ Regm̄ Vni. v? <sup>Ve n̄</sup>  
Ex s̄on̄ Inu. no v? Egredict̄  
In laudib⁹ Antiphona **C** i et do mi rus  
ciuet domini?

et no n t ar d a b i t et ill u m i n a b i t a b s c o n d i  
**et non tardabit et illuminabit abscondi**

X. Ecce te ihesu.  
vnu<sup>m</sup> Splendor, v.

gentibus alleluia ps misere. Emette agnū.

Dicit dominus  
icit dominus peccati a magis  
penitentiam agi Antip.

te appropinquabit enim regnum celorum al-

te appropinquabit cuius regnum celorum al-

leluya Euuae pimā Iherusalem gaudē

leluya Euuae pimā Iherusalem gaudē

# Thank you!