# Disambiguation and Morphological Analysis of Linear B Sign Sequences

# The DĀMOS corpus

- ▶ Only online, (not yet) annotated corpus of Mycenaean Greek
- ▶ Still unlicensed, will be Open Access
- ▶ All 5800 known Linear B documents are available
- ▶ Currently annotation is done manually
  - ▶ Query database for potentially interesting words
  - ▶ Look up rules for sign sequence expansion
  - ▶ Disambiguate
  - ▶ Verify using dictionary
  - ▶ Requires knowledge of ancient Greek, English, and Spanish

# The DĀMOS corpus

- ▶ Write an annotation support tool
- ▶ Let a machine do all the heavy lifting
- ▶ Automatize at much as possible of the pipeline
- ▶ Luckily no syntactic analysis needed (not much syntax in the corpus)

# A Short History of Linear B

- ▶ Script used to write earliest dialect of Greek, Mycenaean
- ▶ Used in Crete and the mainland between 1400 and 1200 BCE
- ▶ Extant texts (mainly on clay tablets) are solely for temporary administrative records
- ▶ Syllabic, i.e. each sign corresponds to a single syllable
- ▶ Some ideograms for objects and units of measure
- ▶ Overall slightly less than 200 unique signs

# A Short History of Linear B

- ▶ Descends from the Linear A script
  - ▶ Predecessor remains undeciphered
- ▶ Current consensus that Linear A was used to write a language unrelated to Greek
- ▶ Decipherment took more than 70 years
- ▶ Groundwork by Alice Kober, full decipherment by Michael Ventris and John Chadwick [Ventris53]

# A Short History of Linear B

- ▶ Construction of combinatoric sign grids
- ▶ Assumption that inflection usually happens in the last syllable and won't change consonant
- ▶ Creation of an internally consistent grammar
- ▶ Guessing a single word and checking if everything else falls into place (a-mi-ni-so → Amnissos)
- ▶ Understanding is still not complete

# A Short History of Linear B



Figure: Syllabary [Ventris53]

# A Short History of Linear B



Figure: Ideograms [Unicode13]

# Methodology

- Syllabary for an unrelated language
- Sign sequences are phonetic approximations of Mycenaean words
- Not able to express all phonetic features of Greek
    - Consonant clusters
    - Ending consonants
    - Diphthongs
- Scribes created various hacks to circumvent script limitations
- Possibly weaknesses in the assigned phonetic values
- Possibly two very similar dialects

# Methodology - Ambiguity and Irregularities

- po-me - ποιμήν - shepherd
  - Dative singular: po-me-ne - ποιμήνέι
  - Nominative singular: po-me-ne - ποιμήνές
- i-qo → ίππος
- i-ko → ιχος or ικος or ισκος
- No syntax to aid reconstruction of intended declension

# Methodology - Preprocessing

- Almost no preprocessing necessary
- Scribes were helpful
  - Comma between words
  - Signs on neatly drawn lines
  - No hyphenation
  - Medium encouraged fixing mistakes
- Considerable effort invested into correctly identifying signs
- Lack of syntax $\rightarrow$ One to one correspondence of lines to sentence-ish constructs

# Methodology - First approaches

- Approximate matching using n-grams
- Create dictionary of Mycenaean words
- Run spell checker on sign sequences
- Produces nonsensical suggestions
  - Phonetic gap too large

# Methodology

- Stochastic language models won't work
- Average text length 14 signs
- Roughly 150 documents exceed 50 signs (not words!)
- Longest text around 600 signs

# Methodology

- Tried-and-True rule based analyzers
- Good track records with other Greek dialects [Crane91]
- Two step process is reduced to single problem
- Easily adapted to state-of-the-art knowledge
    - ...and limited by it

# Methodology

- Drawback: Have to assemble rule set
  - Even worse: there are multiple to choose from
  - Worst of all: all are written in different languages
  - Collaborator agreed to assist

# Methodology - Basic Script Rules

- $p \to p^h, b, p$
- $k \to k^h, g, k$
- $r \to r, l$
- $q \to k^w, g^w$
- Syllables with initial consonant clusters are written with signs matching the vowel of the syllable
- Ending consonants are omitted (ti-ri-po-de $\to$ tripodes)
- Much more complex ones possible (some site specific [Hooker80])

# Methodology - Measurement

- Measurement of success/failure quite complex
- No training set to verify against
- Did build Ad-Hoc references from literature for ngram similarity test
- Have to rely on subjective impressions from collaborators

# Outlook

- Assemble rule set
- Convert it into a sensible format
  - Using endless chains of regular expressions is not satisfactory
  - In an ideal world rule set and source code is kept apart
- Functional Morphology sounds appropriate

# Outlook

- Plan to have robust results until finals
- Maybe incorporate the ngram similarity algorithm as a postprocessing step
- Try to extend tool's application to more than annotation assistance
  - If output is stable across known similar texts look for notably worse results
- Unfortunately, irreducibly complex

# Propaganda

- Bring civili...computer science to the barbarians
- It hurts to see how they do things
- Find/Show appropriate tools for working with small and barely understood information systems
- Hopefully, won't need constant funding (and maintenance) to be of use

# References

📄 Evidence for Greek Dialect in the Mycenaean Archives, Ventris and Chadwick, The Journal of Hellenic Studies, Vol. 73 (1953), pp. 84-103

📄 Generating and parsing classical Greek, Crane, Literary and Linguistic Computing 6.4 (1991) pp. 243-245

📄 The Unicode Standard 6.3

📄 Linear B: An Introduction, Hooker, Bristol Classical Press, 1980