

Sunoikisis Digital Classics
Fall 2020: session 4

Overview of Computational Linguistics

Alek Keersmaekers (KU Leuven)
Marton Ribary (Surrey)
Thea Sommerschield (Oxford)

Classical Computational Linguistics: an introduction

What is computational linguistics?

- In a broad sense: any linguistic work with a computational methodology
- Computational analysis of linguistic data (i.e. corpus linguistics)
- Automated processing and/or annotation of linguistic data (i.e. natural language processing)



computational linguistics

noun

the branch of linguistics in which the techniques of computer science are applied to the analysis and synthesis of language and speech.

Corpus compilation (1)

- Most of the work in classical computational linguistics: focused on the compilation of corpus resources
- Full-text databases (see also https://wiki.digitalclassicist.org/Greek_and_Latin_texts_in_digital_form):
 - *TLG, Perseus* (literary texts)
 - *Packard Humanities Institute* (inscriptions)
 - *Duke Databank of Documentary Papyri* (papyri)
- Automatically lemmatized + morphologically annotated texts: e.g. *Perseus under Philologic, Diorisis, LatinISE*
- Manually lemmatized + morphologically and syntactically annotated texts (treebanks): see next slide
- Has led to several corpus-based studies (even though the classical languages are still lagging behind: see Jenset & McGillivray 2017)

Corpus compilation (2)

Greek project	Tokens	Texts
AGDT	560K	Archaic + Classical poetry/prose
Gorman	324K	Classical + post-classical prose
Pedalion	300K	Classical + post-classical poetry/prose
PROIEL	270K	Herodotus, NT, Sphrantzes
Harrington	18K	Lucian, Septuagint, Life of Aesop
Aphthonius	7K	Aphthonius
Sematia	6K	Papyri

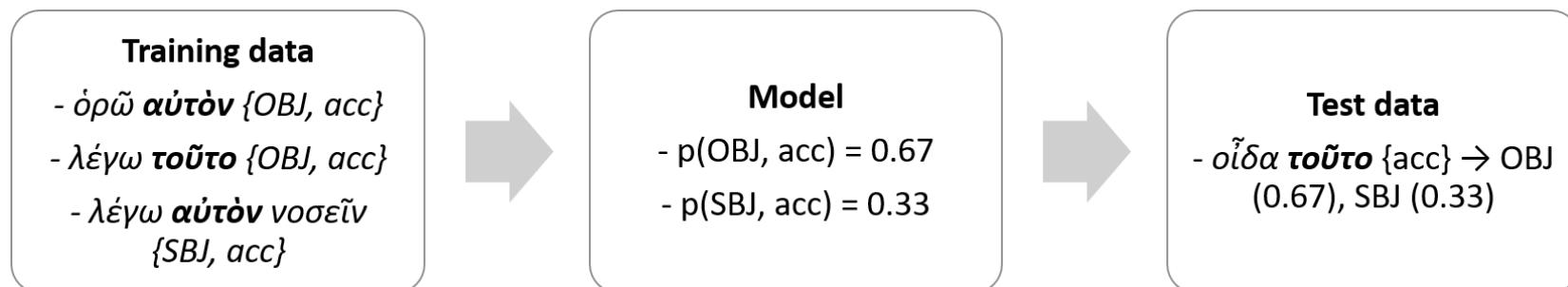
Latin project	Tokens	Texts
Index Thomisticus	450K	Thomas Aquinas
PROIEL	225K	Classical + post-classical prose
LLCT	200K	Medieval charter texts
Harrington	120K	Classical prose/poetry
ALDT	53K	Classical prose/poetry

NLP and Classical languages

- ... in contrast, natural language processing of classical languages is still in its infancy (large amounts of data required)
- Natural language can be processed in several ways:
 - Morphological processing and lemmatization (McGillivray & Kilgariff 2013; Celano, Crane & Majidi 2016)
 - Syntactic parsing (e.g. Mambrini & Passarotti 2012; Ponti & Passarotti 2016)
 - Named Entity Recognition (e.g. Erdmann et al. 2016; Palladino, Karimi & Mathiak 2020)
 - Distributional lexical semantics (e.g. Rodda, Senaldi, and Lenci 2016; Bamman and Burns 2020)
 - ...
- See the 3 case studies in this session for more examples
- See also the Classical Language Toolkit (<http://cltk.org/>, see also [https://github.com/SunoikisisDC/SunoikisisDC-2017-2018/wiki/The-Classical-Language-Toolkit-\(CLTK\)](https://github.com/SunoikisisDC/SunoikisisDC-2017-2018/wiki/The-Classical-Language-Toolkit-(CLTK)) and <https://www.digitalclassicist.org/wip/wip2018-09pb.html>)

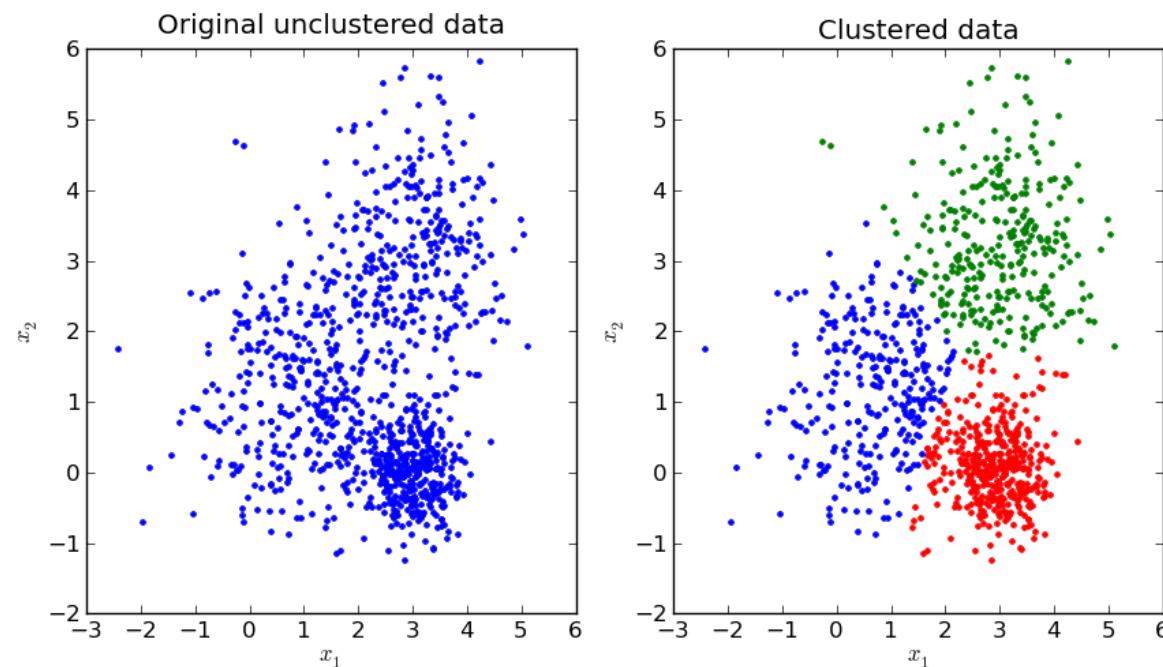
Machine learning: main concepts (1)

- Older computational work (for English): rule-based
- Most methods nowadays: machine learning
- **Supervised** machine learning:
 - To predict a certain class (e.g. OBJ vs. SBJ)
 - ... one starts from a large corpus annotated for a number of **features** (e.g. case), the **training data**
 - The computer ‘learns’ patterns from the data and formalizes them in a mathematical **model**
 - This model can be used to predict class labels for new, unseen data, the **test data**



Machine learning: main concepts (2)

- **Unsupervised** machine learning (e.g. clustering): no pre-defined class labels
- Instead, the computer divides the data into groups with internally similar members



Classical computational linguistics: main challenges

- Typological characteristics of Greek and Latin (most NLP methods tailored to English): e.g. highly inflectional, free word order, large rate of ellipsis
- Large genre and diachronic variation
- Data sparseness

References

- Bamman, D. & Burns, P. J. 2020. Latin BERT: A Contextual Language Model for Classical Philology. *arXiv Preprint arXiv:2009.10053*.
- Celano, G. G. A., Crane, G. & Majidi, S. 2016. Part of Speech Tagging for Ancient Greek. *Open Linguistics* 2(1).
- Erdmann, A., Brown, C., Joseph, B., ... Marneffe, M.-C. de. 2016. Challenges and Solutions for Latin Named Entity Recognition. In *Proceedings of the Workshop on Language Technology Resources and Tools for Digital Humanities (LT4DH)*. Osaka: The COLING 2016 Organizing Committee, pp. 85–93.
- Jerset, G. B. & McGillivray, B. 2017. *Quantitative historical linguistics: a corpus framework*. Oxford: Oxford University Press.
- Mambrini, F. & Passarotti, M. C. 2012. Will a parser overtake Achilles? First experiments on parsing the ancient Greek dependency treebank. In *Eleventh International Workshop on Treebanks and Linguistic Theories*. Edições Colibri, pp. 133–144.
- McGillivray, B. & Kilgarriff, A. 2013. Tools for historical corpus research, and a corpus of Latin. In Durrell, P., Scheible, M., Whitt, S. & Bennett, R. J. (eds.), *New Methods in Historical Corpus Linguistics*. pp. 247–257.
- Palladino, C., Karimi, F. & Mathiak, B. 2020. NER on Ancient Greek with minimal annotation. In *Digital Humanities 2020*. Ottawa: DH2020.
- Ponti, E. M. & Passarotti, M. 2016. Differentia compositionem facit. A Slower-Paced and Reliable Parser for Latin. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*. Portorož: European Language Resources Association (ELRA), pp. 683–688.
- Rodda, M. A., Senaldi, M. S. & Lenci, A. 2016. Panta Rei: Tracking Semantic Change with Distributional Semantics in Ancient Greek. In *CLiC-it/EVALITA*.

Text pre-processing and semantic clustering

Marton Ribary (Surrey)

Sunoikisis Digital Classics
Fall 2020 Session 4. Computational Linguistics
Thursday 29 October 2020

Article

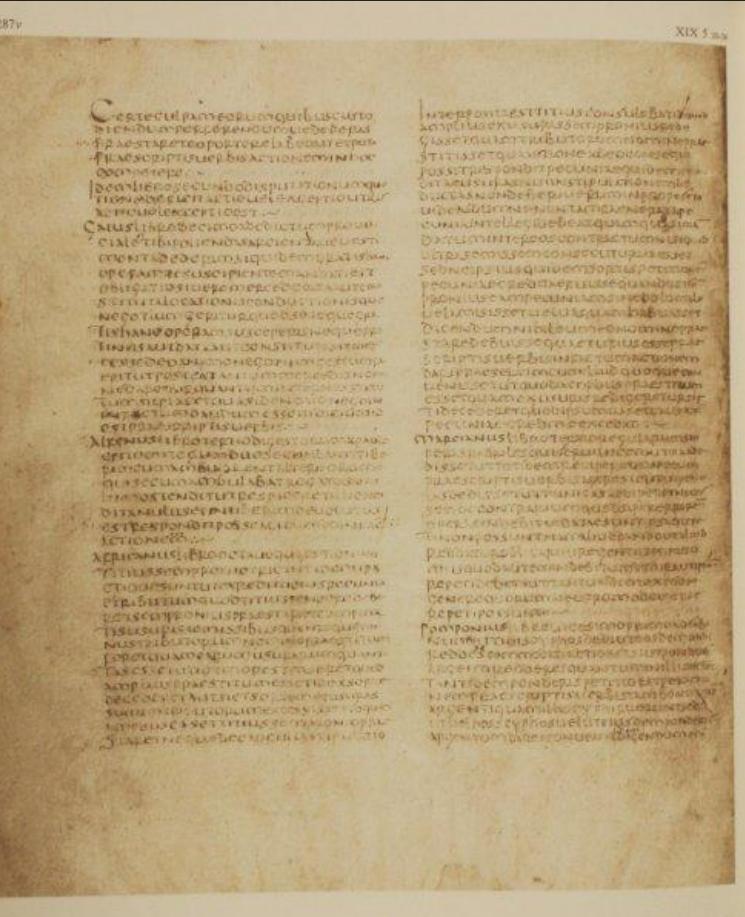


Ribary, M and McGillivray, B.,
“A Corpus Approach to Roman Law Based on
Justinian’s Digest.”
Informatics 7 (2020), 44.
DOI: <https://doi.org/10.3390/informatics7040044>

Demo



<https://colab.research.google.com/drive/1fBsLYLeuh4raQ1CGfvIWAZe2ftKRzHXR?usp=sharing>



Manuscript

- Digest - historical sourcebook of Roman law in 432 thematic sections
 - 9132 passages in 21055 text units
 - (almost) exclusive source for Roman law
- compiled by Tribonian under the order of Justinian I
- and published in 533 CE
- oldest (and complete) MS from 554 CE

dem tibi inspicendam dedi et dicas te perdidisse, ut demum mili praeceptis tuis, quae incepimus, si memoremur, tibi agere possum vel conducere vel ad exhibendum agere. Secundum haec, si cui inspicendum dedi sive ipsius causa utrrixque, et datus et cupimus mili praeceptis tuis, quae incepimus, tibi agere possum; si vero me duxeratis, causa datum est, dolum 3 solum, quia prope depositum habeo accessit. Si, cum unum bovem habrem et vicinus unum, plausim mili bovem commendares, ut opus faceret, et apud alterum bovis perit, commodati non competit actio, sed iure libro secundo disputationis Quoties dictum acutum est, *quod deinde a dolo*.

21. *Iure libro secundo disputationis Quoties dictum acutum est, quod deinde a dolo*

tibi polliuuntur sacerdotes vestimenta sacerdotum, si mandat gratis hanc operam a susceptore, mandati est obligatio, si vero mercede aut acto constituta, loca-
tio etiam obligatio, et si mandat gratis hanc operam susceptori neque pretiosum aut data aut constituta sit merces, sed eo nomine no-
getum gestum fuerit, at post tantum mercede no-
getum gestum fuerit, statim et statim sit, placet
quid est novo negotio in verbis factum dandum esse hui-
cium quod est praeceptum in verbis factum dandum esse hui-

33 ALFENUS tertio digestorum a Paulo epitolam Dux secundum Tiberium cum ambularet, alter eum ei, qui secum ambulabat, rogatus annui ostendit, ut responderet: illi excedit annus et Tiberius devolutus est. respondit posse agi cum eo in factum actione.

15 *Tinere.* Non sibi sponsonis causa ausus accepterit ne reddit² vitori, prescripcis verbis actis in eum competit ne amissio recipienda est Sabini opinio. **16** *Si tamen et furi ait.* Haec causa non permissum modum enim ostendit, quae sponsonis plane neque dominica vitor habuit, agit furi? plane si in honesta causa sponsonis fulit, sui furi? plane dumtaxat repetito erit.

pecuniam depositurum, ut dares Titio, si fugitum
meum redisset, nec deudis, quia non redixit: si
pecunia mibi non reddita, melius est prescriptis
verbis agere: non enim alio pecuniam ego et fugi-
tivarum depositum, ut quasi apud sequestrem sit
depositum.

19 Item libro trigesimo pro aliud ad eundem Rogasti me, ubi tribus mutatos daret: eum qui non haberet, dili tibi rem vendendum, ut per te uteris, si non vendidisti aut vendidisti quidem, pecuniam tuam, et non videris nisi videris, et non agas cum aliis. Labo ad praescripsit verbis, quasi nego, ita quodam inter nos gesto propterea contraria. Si prae nos, ut mihi fidei obligavero, deinde placuerit inter nos, nam propter praestares, ne facias, metuens esse dico praescripsit verbis agi, nisi merecere interveniat: ne si intervenerit verbis loco esse actionem.

20 IDEM LIBRO triginto secundo ad edictum Apud Labecum queritur, si tibi quis venales experientur decessus at, si in triadis displicenter reperiatur, quod non possit et videretur decessus emere nolueris, si ad adversus te ex venditu actio. et puto verus esse prescripsit verbis agendum: nam inter nos hoc actum, ut experimentum gratuitum accepimus, non est enim cetero quod possit. sed si mulier te decessu te exprimere ac placuisse, scilicet si multo te decessu te exprimere ac placuisse, scilicet si displicenter et in dies singulos aliquid praestares, deinde mulier a grossorum frumenti auribus alata intra dies experimentum quod esset praestare, et non prestatum, et quod tunc erit iustitia, et alii dicunt, ut utrumque, utrum empicio iam erat contracta ad futurum, ita, ut si^t facta, preludetur, si futura, merces petatur: sed non exprimit actionibus, puto autem, ut expeditius, praefacta fuit emptio contracta ex venditu, et si dicitur, quod taliter taliter petatur, etiam, quoniam taliter aduersus 2. desuperum est, scilicet si cum emere argenteum velles, vasularium at de detulerit et reliquerit, et cum dicas pliecessit tibi, servu tuo referendum deducit et sine placuisse et cum tibi detulerit et reliquerit, et cum dicas meum, quod taliter causa sit missam, certe culpm eorum, quibus custodiendum perferendum defers, praestare te oportet. Labet si, et puto prescripsit verbis actionem in hunc competere.

21 IDEM LIBRO triginto secundo ad edictum Apud Labecum queritur, si tibi quis venales experientur decessus at, si in triadis displicenter reperiatur, quod non possit et videretur decessus emere nolueris, si ad adversus te ex venditu actio. et puto verus esse prescripsit verbis agendum: nam inter nos hoc actum, ut experimentum gratuitum accepimus, non est enim cetero quod possit. sed si mulier te decessu te exprimere ac placuisse, scilicet si multo te decessu te exprimere ac placuisse, scilicet si displicenter et in dies singulos aliquid praestares, deinde mulier a grossorum frumenti auribus alata intra dies experimentum quod esset praestare, et non prestatum, et quod tunc erit iustitia, et alii dicunt, ut utrumque, utrum empicio iam erat contracta ad futurum, ita, ut si^t facta, preludetur, si futura, merces petatur: sed non exprimit actionibus, puto autem, ut expeditius, praefacta fuit emptio contracta ex venditu, et si dicitur, quod taliter taliter petatur, etiam, quoniam taliter aduersus 2. desuperum est, scilicet si cum emere argenteum velles, vasularium at de detulerit et reliquerit, et cum dicas pliecessit tibi, servu tuo referendum deducit et sine placuisse et cum tibi detulerit et reliquerit, et cum dicas meum, quod taliter causa sit missam, certe culpm eorum, quibus custodiendum perferendum defers, praestare te oportet. Labet si, et puto prescripsit verbis actionem in hunc competere.

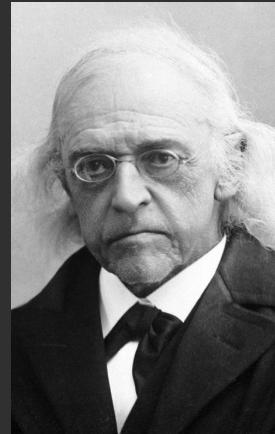
22 MARECIUS LIBRO tertio regularum Si oras fabriles possunt esse vice mensis, et tunc tamen si pauculas possunt, cum variis signis accepit: nec esse habeat contrarium, quid, si per errorem operae inadhibe date contraria, quis ipsa reged possumit nam aliud date, et si per errorem operae inadhibe date, quis ipsa reged possumit, etiam si per errorem operae inadhibe date, prius potest autem inadhibe date, ipsius neuro modo opera reged possunt.

23 POMONIUS LIBRO vicecentesimo et ad Solisiam Si tunc est omnes dicti, ut comedimus menses, et componimus isti, et si, quod tamponis ardentis reddimus cum illis easel, tamponis ardenti petito est per actionem prescripsit verbis, tam boni tamen argenti, quam illi sycophantibus: sed si ut vel hoi sycophantibus est, et si quidam qui sunt electi, et sycophantibus est, et si quidam qui sunt electi, sycophantibus est, tunc huius feri et te mibi dare aut sycophantibus argenteum utrum malis: quod si nulli mihi sumptibus est eligere, sycophantibus non fieri antea diximus, et si nulli mihi sumptibus est eligere, sycophantibus non fieri, et si doli solidis medicis et scientias, non est merces: itaque si quid in his ministeris alter fiat quam convenient, non ex locato, sed in factum actio da-
bitur.¹¹

(1) qui F (2) vieti ins. secundum B (3) mulam F^2
 (4) et merces del. (5) sit ins. (6) sarvientiae
 (7) vaquana F (8) quae F (9) qui F (10) reddes F | (11) ante dicendum est ins. idem F^2 : quae sequuntur omissa
 in F supplentur secundum B

Printed edition

- Largely based on the *littera Florentina*
 - Edited by Theodor Mommsen



Theodor Mommsen (1817-1903)

Digital edition

Amanuensis V5.0

Search About

Enter your search terms Ελληνικά γράμματα Stay on top

d.19,5,

Start searching

brought to you by
ERC Starting Grant >ACO<

erc

Search results Context ?

D. 19, 5, 17, 3 Ulp. 28 ad ed.
Si, cum unum bovem haberem et vicinus unus, placuerit inter nos, ut per denos dies ego ei et ille mihi bovem commodaremus, ut opus faceret, et apud alterum bos periit, commodati non competit actio, quia non fuit gratuitum commodatum, verum praescriptis verbis agendum est.

D. 19, 5, 17, 4 Ulp. 28 ad ed.
Si, cum mihi vestimenta venderes, rogavero, ut ea apud me relinquas, ut peritoribus ostenderem, mox haec perierint vi ignis aut alia maiore, periculum me minime praestaturum: ex quo apparet utique custodiam ad me pertinere.

D. 19, 5, 17, 5 Ulp. 28 ad ed.
Si quis sponzionis causa anulos acceperit nec reddit victori, praescriptis verbis actio in eum competit: nec enim recipienda est Sabini opinio, qui condici et furti agi ex hac causa putat: quemadmodum enim rei nomine, cuius neque possessionem neque dominium vicitur habuit, aget furti? plane si inhonesta causa sponzionis fuit, si anuli dumtaxat repetitio erit.

D. 19, 5, 18, 0 Ulp. 30 ad ed.
Si apud te pecuniam depositum, ut dares Titio, si fugitivum meum reduxisset, nec dederis, quia non reduxit: si pecuniam mihi non reddas, melius est praescriptis verbis agere: non enim ambo pecuniam ego et fugitivarius depositimus, ut quasi apud sequestrem sit depositum.

D. 19, 5, 19, 0 Ulp. 31 ad ed.
Rogasti me, ut tibi nummos mutuos darem: ego cum non haberem, dedi tibi rem vendendam, ut pretio uteroris. si non vendidisti aut vendidisti quidem, pecuniam autem non acceperisti mutuam, tutius est ita agere, ut Labeo ait, praescriptis verbis, quasi negotio quodam inter nos gesto proprii contractus.

D. 19, 5, 19, 1 Ulp. 31 ad ed.
Si oraedium oro te oblioquero. deinde placuerit inter nos. ut mihi

© 2020 Günther Rosenbaum & Peter Riedlberger

- ROMTEXT database (Linz) since the 1970s
- Amanuensis interface developed by Peter Riedlberger

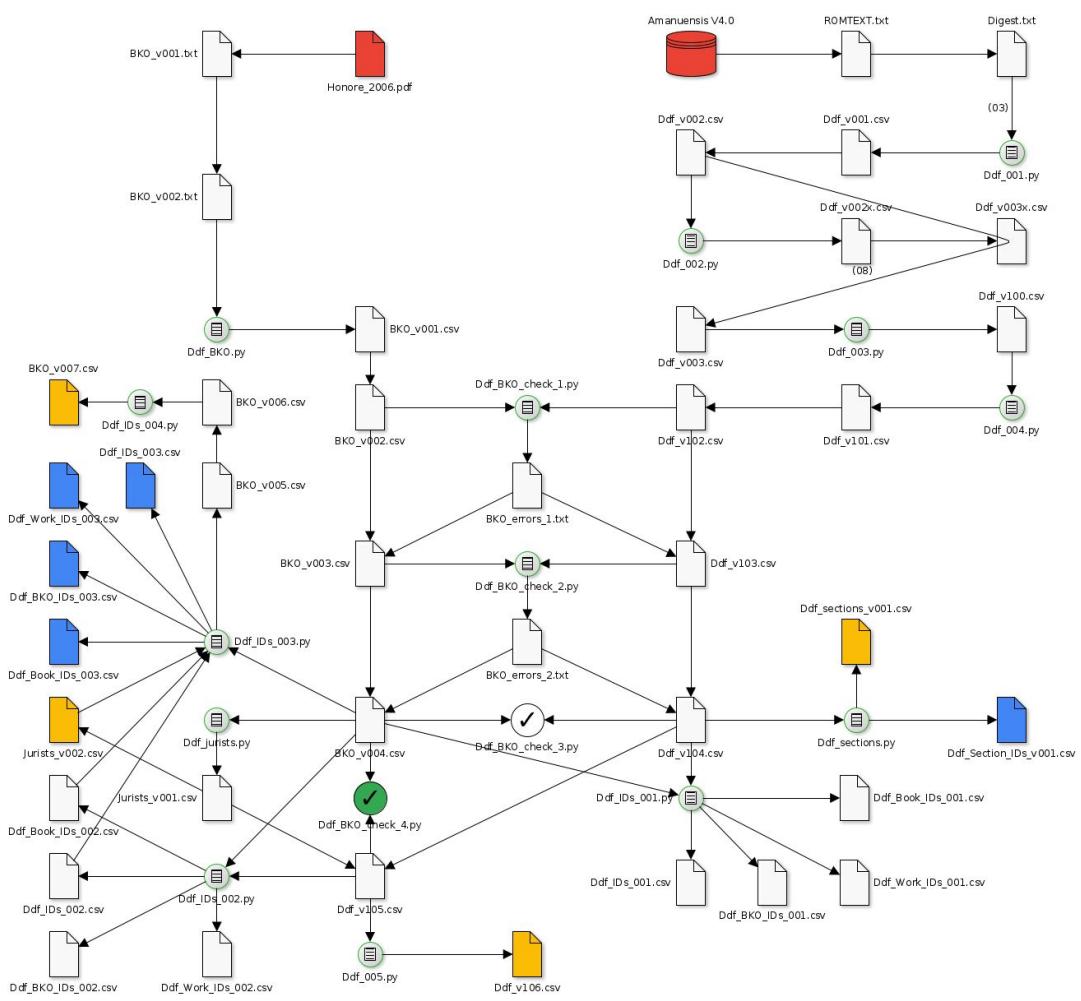
Database



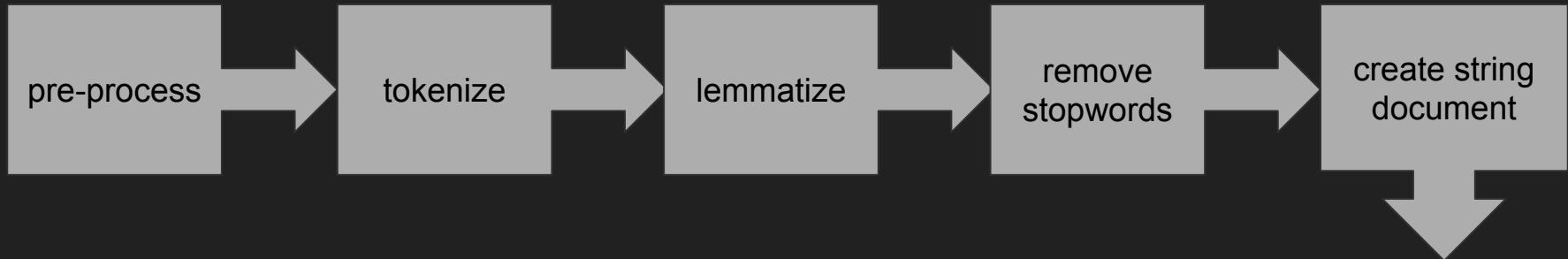
Processing flowchart from the GitLab repository
<https://gitlab.eps.surrey.ac.uk/mr0048/pydigest>

Ribary, M., “A Relational Database of Roman Law Based on Justinian’s Digest.” *Journal of Open Humanities Data* 6(1), p.5. DOI:
<http://doi.org/10.5334/johd.17>

Ribary, M., “A relational database of Roman law based on Justinian’s Digest.” *figshare*. Dataset.
<https://doi.org/10.6084/m9.figshare.12333290.v1>



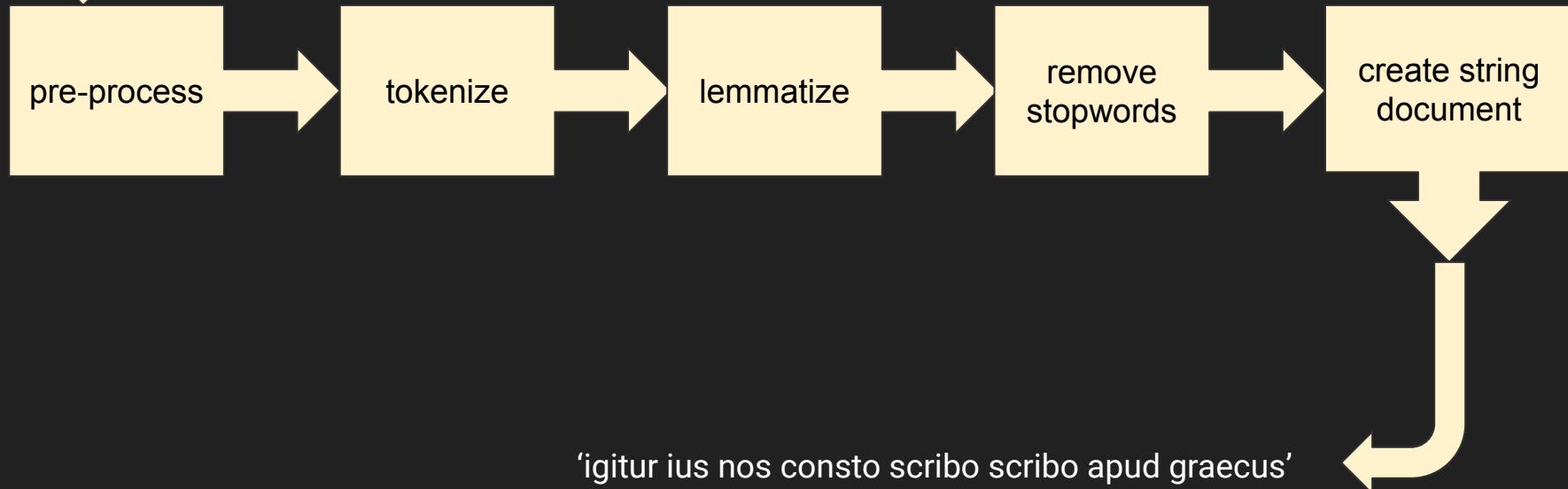
Getting the text ready



[Demo available](#)

'Hoc igitur ius nostrum constat aut ex scripto aut sine scripto, ut apud Graecos: τῶν νόμων οἱ μὲν ἔγγραφοι, οἱ δὲ ἄγγραφοι.'

(D.1.1.6.1)

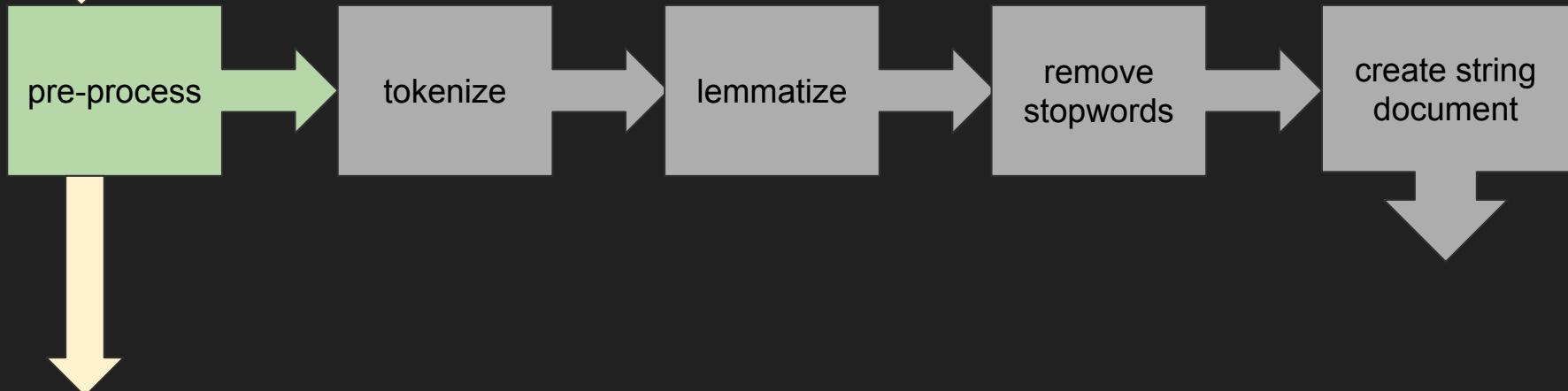


'igitur ius nos consto scribo scribo apud graecus'

[Demo available](#)

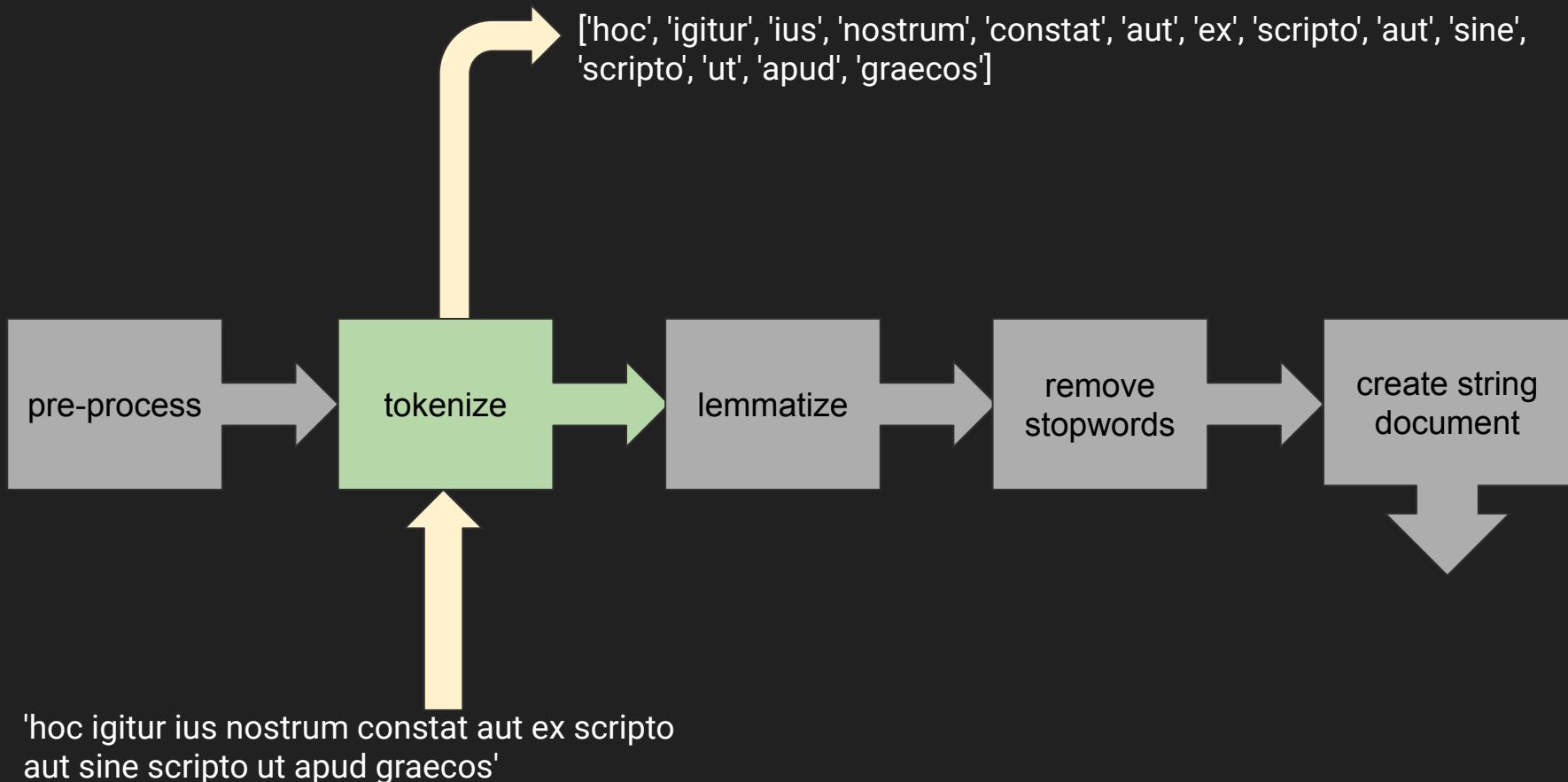
'Hoc igitur ius nostrum constat aut ex scripto aut sine scripto, ut apud Graecos: τῶν νόμων οἱ μὲν ἔγγραφοι, οἱ δὲ ἄγραφοι.'

(D.1.1.6.1)



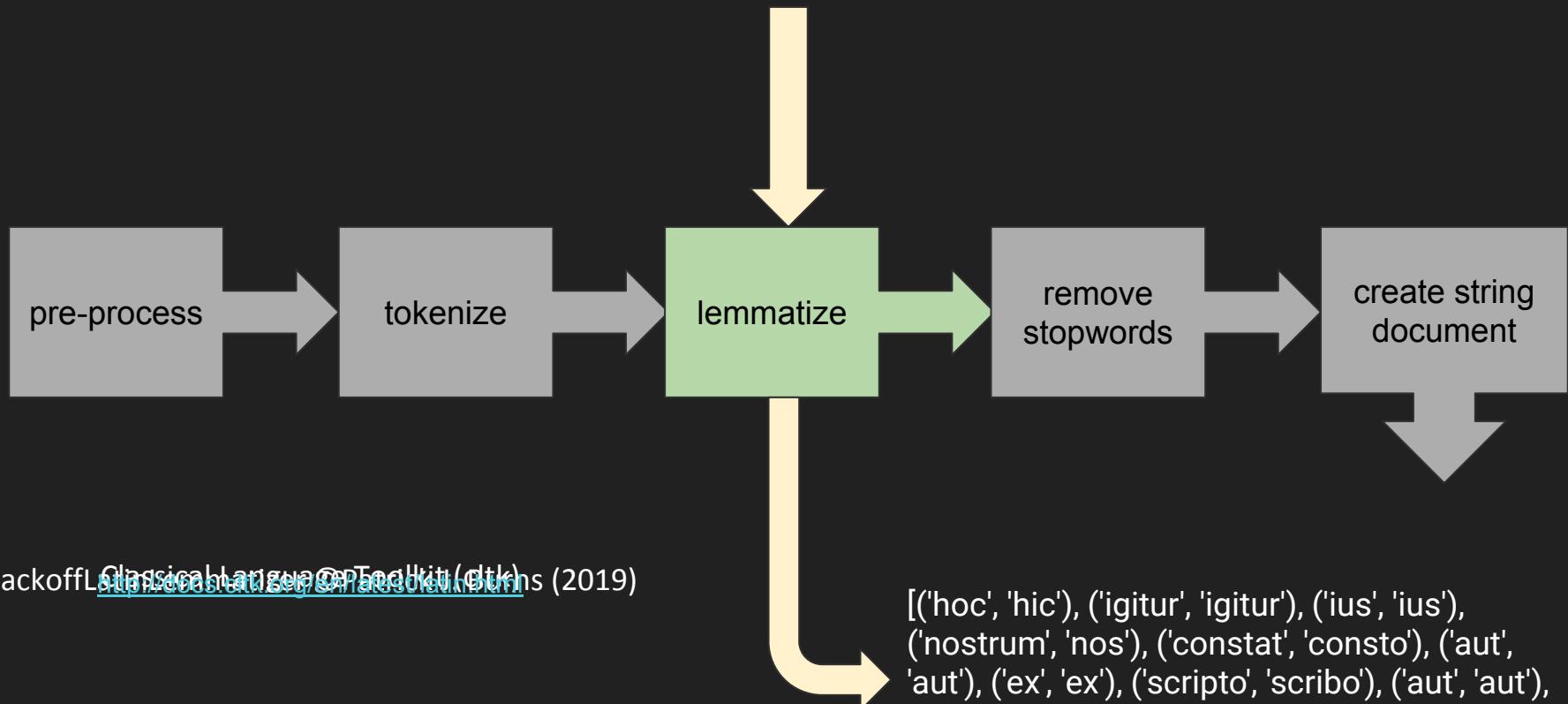
'hoc igitur ius nostrum constat aut ex scripto
aut sine scripto ut apud graecos'

[Demo available](#)



[Demo available](#)

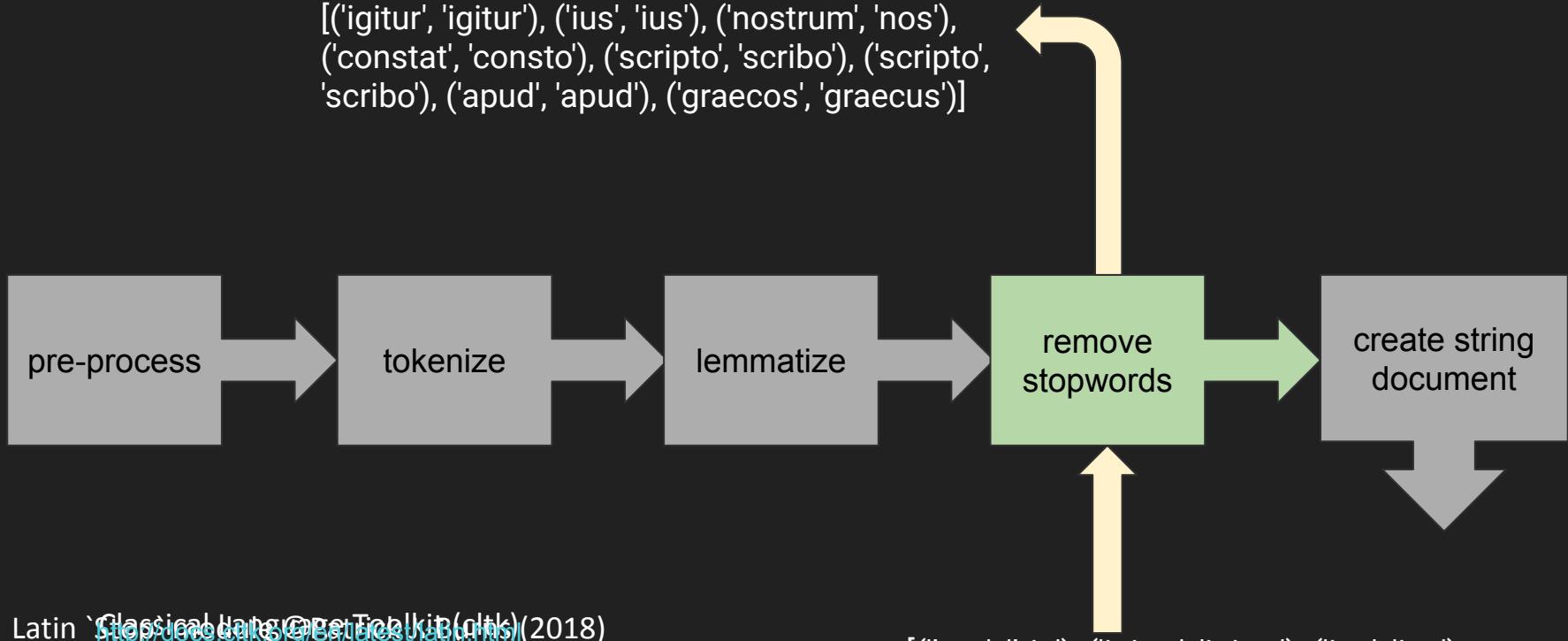
```
['hoc', 'igitur', 'ius', 'nostrum', 'constat', 'aut', 'ex', 'scripto', 'aut', 'sine',  
'scripto', 'ut', 'apud', 'graecos']
```



BackoffLCL Classical Language Toolkit (ctk) (2019)

<https://classicallanguagekit.github.io/>

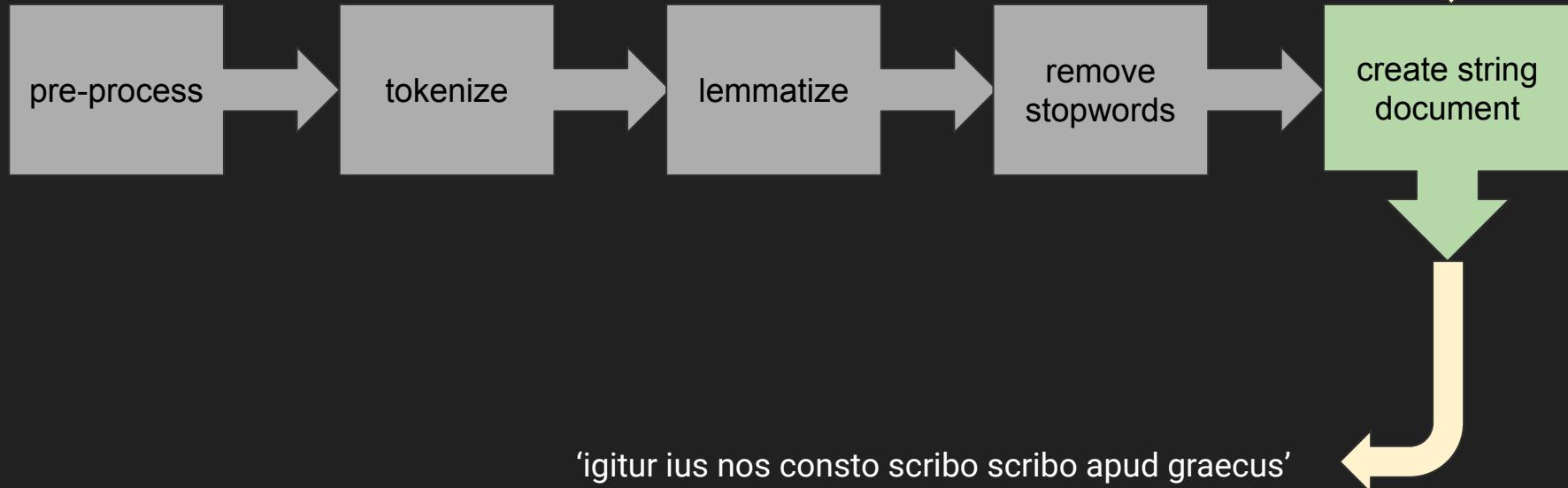
Demo available



[('hoc', 'hic'), ('igitur', 'igitur'), ('ius', 'ius'),
('nostrum', 'nos'), ('constat', 'consto'), ('aut',
'aut'), ('ex', 'ex'), ('scripto', 'scribo'), ('aut', 'aut'),
('sine', 'sine'), ('scripto', 'scribo'), ('ut', 'ut'),
('apud', 'apud'), ('graecos', 'graecus')]

[Demo available](#)

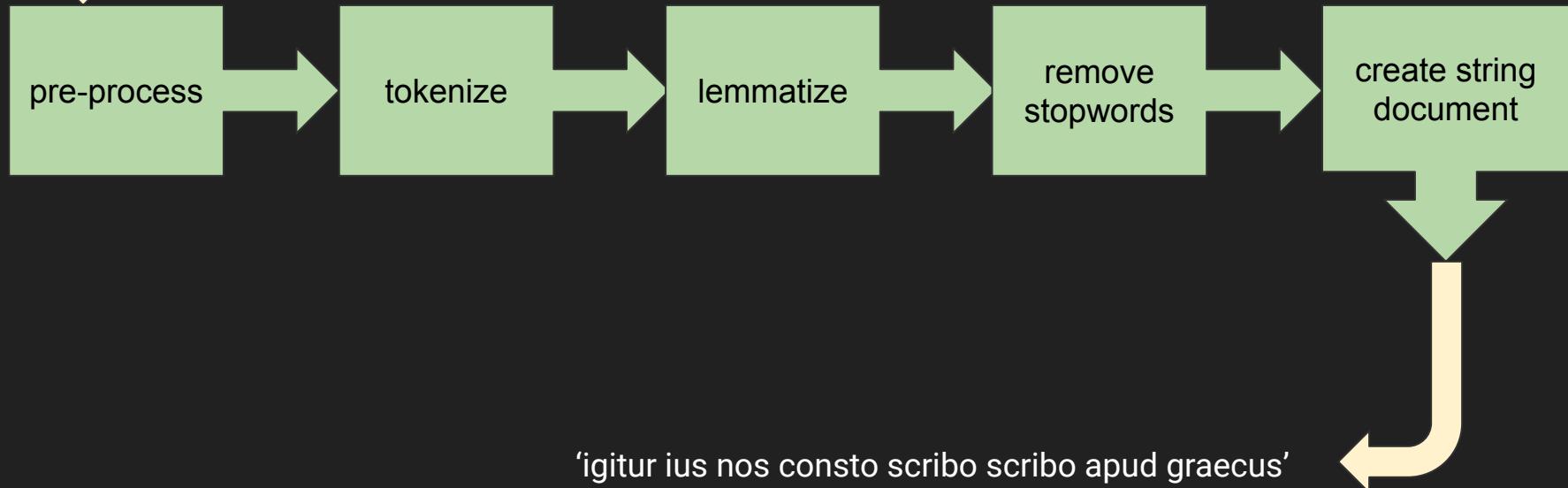
```
[('igitur', 'igitur'), ('ius', 'ius'), ('nostrum', 'nos'),  
('constat', 'consto'), ('scripto', 'scribo'), ('scripto',  
'scribo'), ('apud', 'apud'), ('graecos', 'graecus')]
```



[Demo available](#)

'Hoc igitur ius nostrum constat aut ex scripto aut sine scripto, ut apud Graecos: τῶν νόμων οἱ μὲν ἔγγραφοι, οἱ δὲ ἄγγραφοι.'

(D.1.1.6.1)



[Demo available](#)

Discover the structure of
Roman law from the bottom up



Clustering

Assumptions:

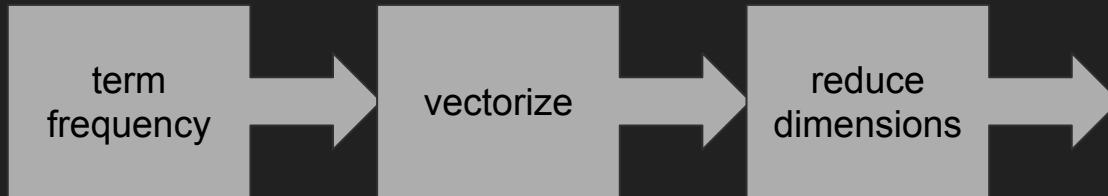
1. Sections are thematically homogeneous
2. The Digest is a comprehensive overview of what Roman law is

Textual data of sections can be transformed to numbers by tfidf

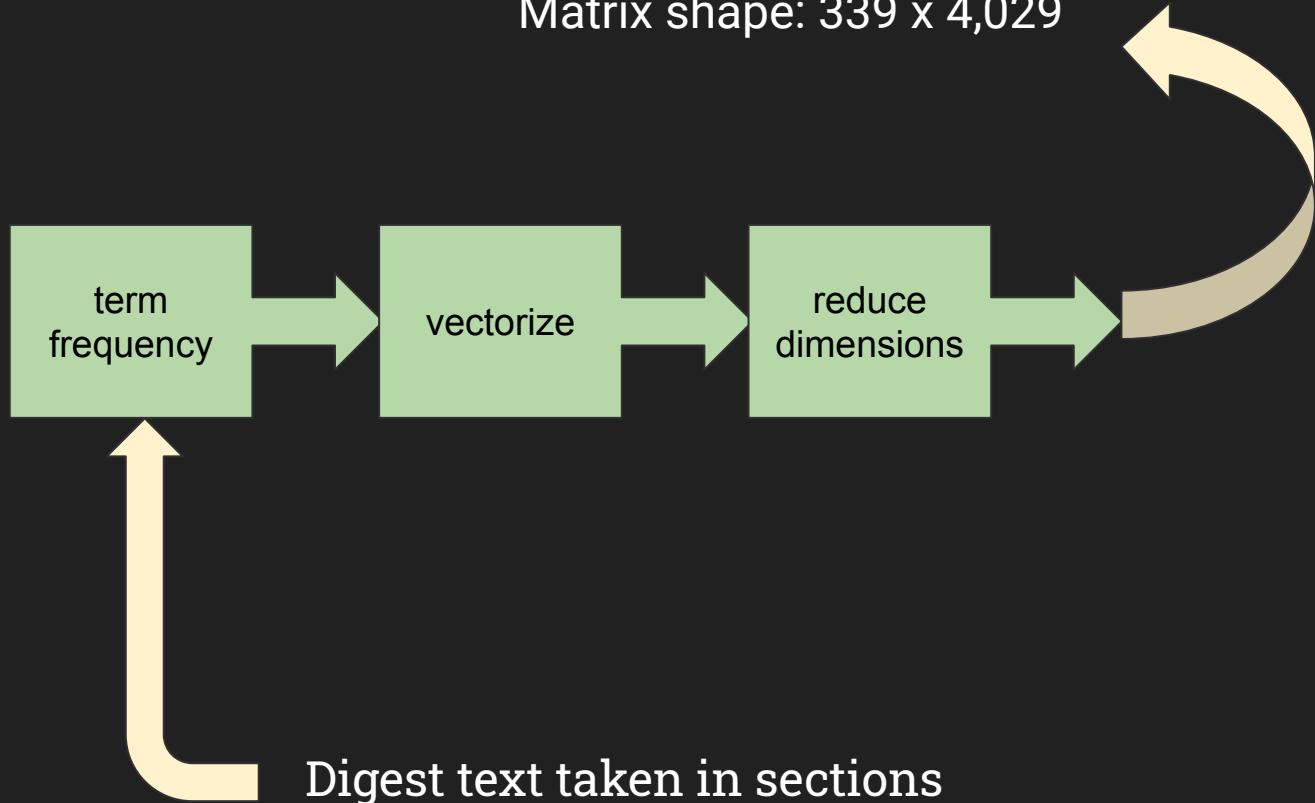
Hierarchical clustering of sections brings out the bottom-up structure of Roman law

Turn text to numbers

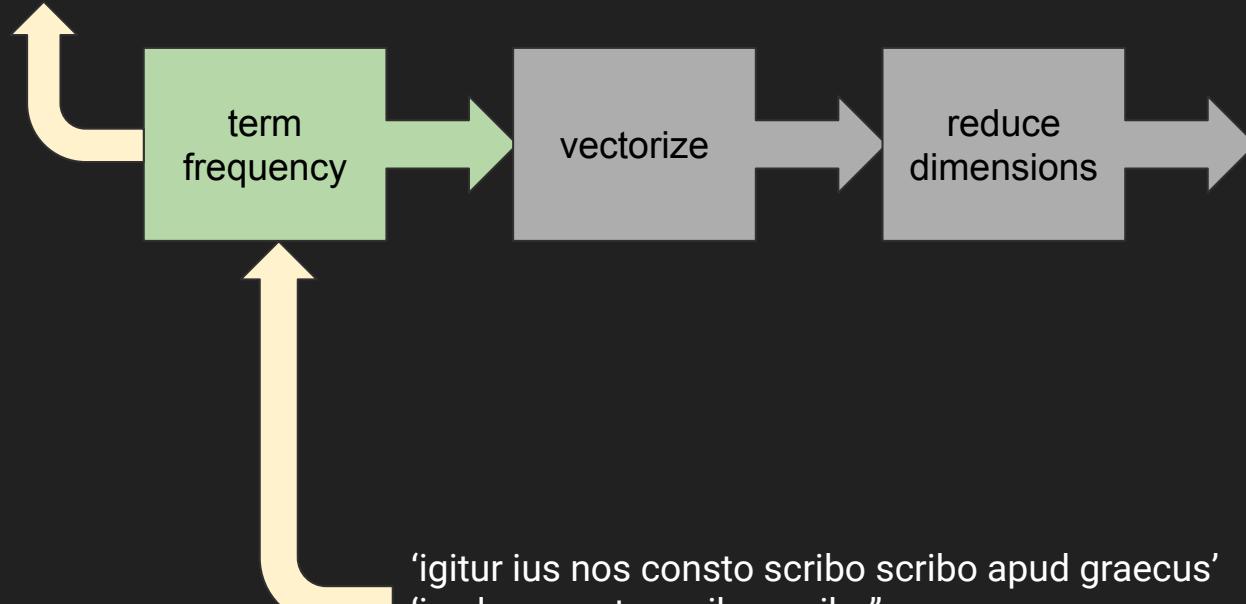
Term frequency - Inverse document frequency (tfidf)



Matrix shape: 339 x 4,029

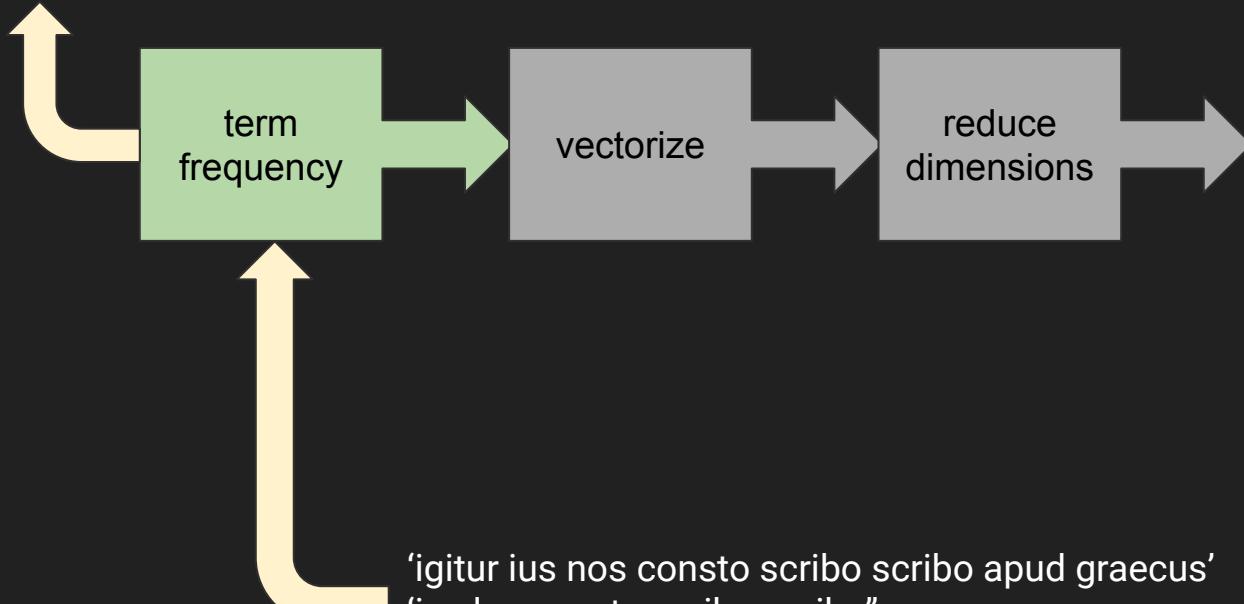


	igitur	ius	nos	consto	scribo	apud	graecus	lex	romanus
1	1	1	1	1	2	1	1	0	0
2	0	1	0	1	2	0	0	1	0
3	1	0	0	1	0	1	0	1	1



'igitur ius nos consto scribo scribo apud graecus'
 'ius lex consto scribo scribo'
 'igitur lex consto apud romanus'

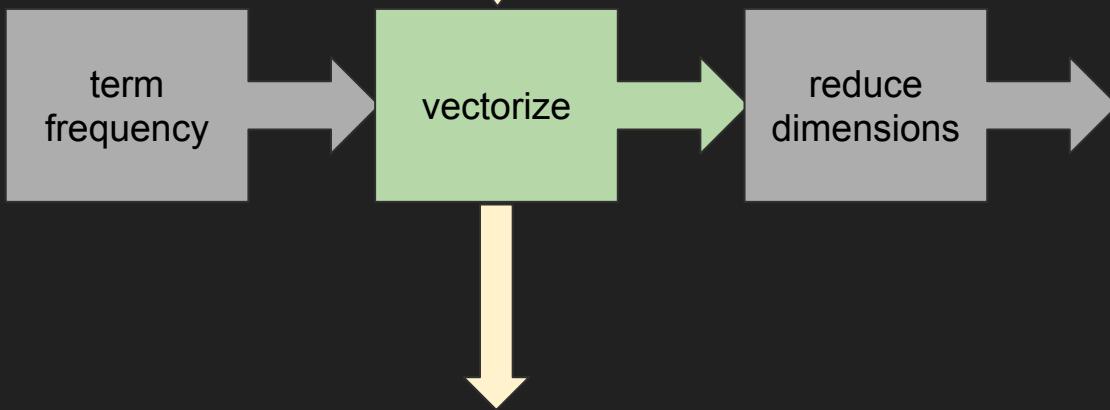
1	1	1	1	2	1	1	0	0
0	1	0	1	2	0	0	1	0
1	0	0	1	0	1	0	1	1



'igitur ius nos consto scribo scribo apud graecus'
'ius lex consto scribo scribo'"
'igitur lex consto apud romanus'

1	1	1	1	2	1	1	0	0
0	1	0	1	2	0	0	1	0
1	0	0	1	0	1	0	1	1

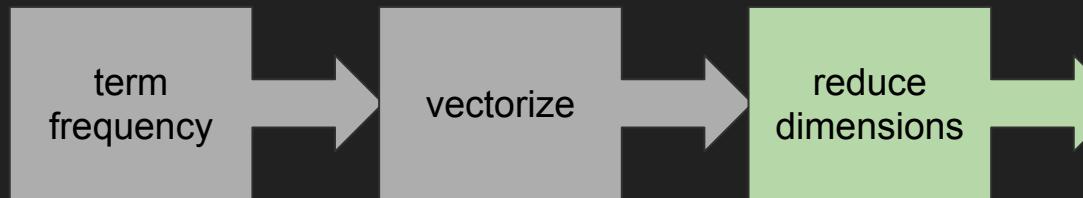
Term frequency - Inverse
document frequency (tfidf)
all sections with word
weighting



Matrix shape: 432 x 10,865
(432 sections, 10,865 terms)

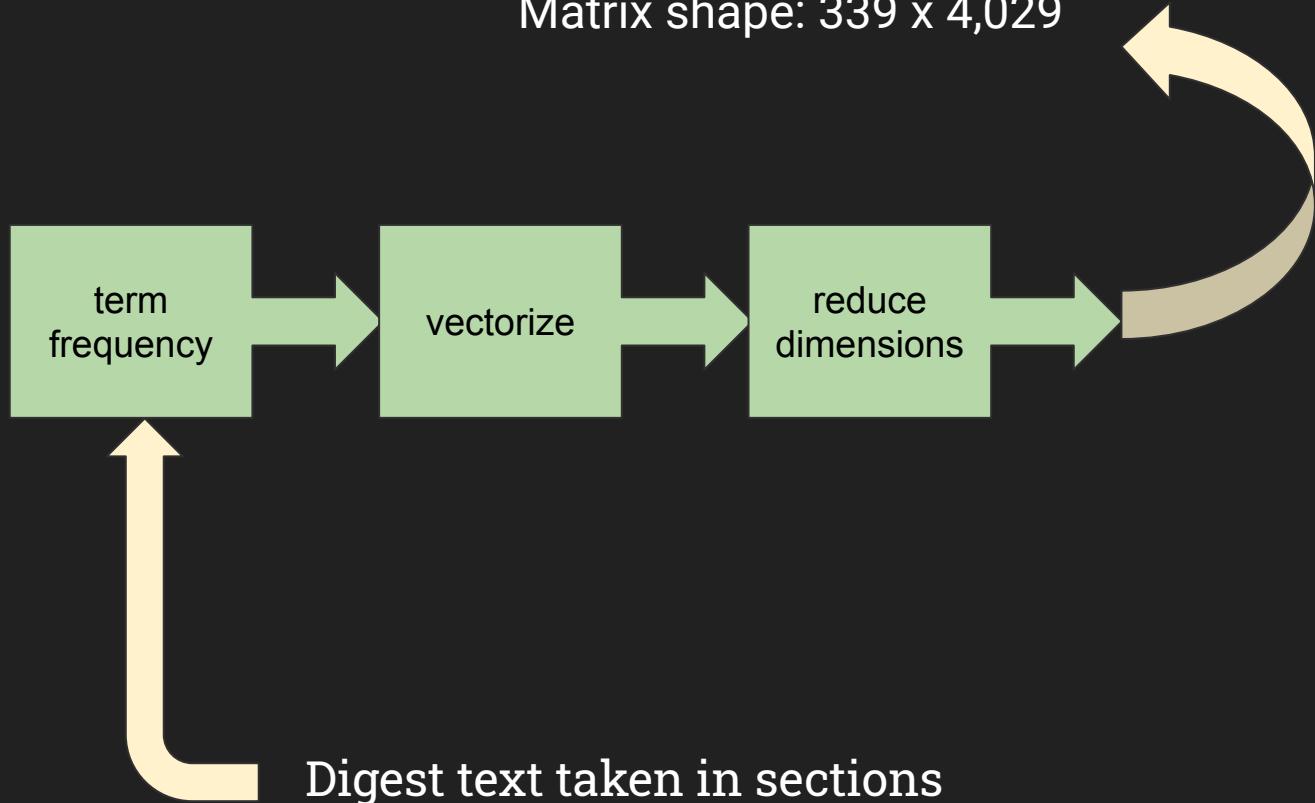
- Remove short sections
(339 remains)
- Select top 50 terms only

Matrix shape: 339 x 4,029
(339 sections, 4,029 terms)

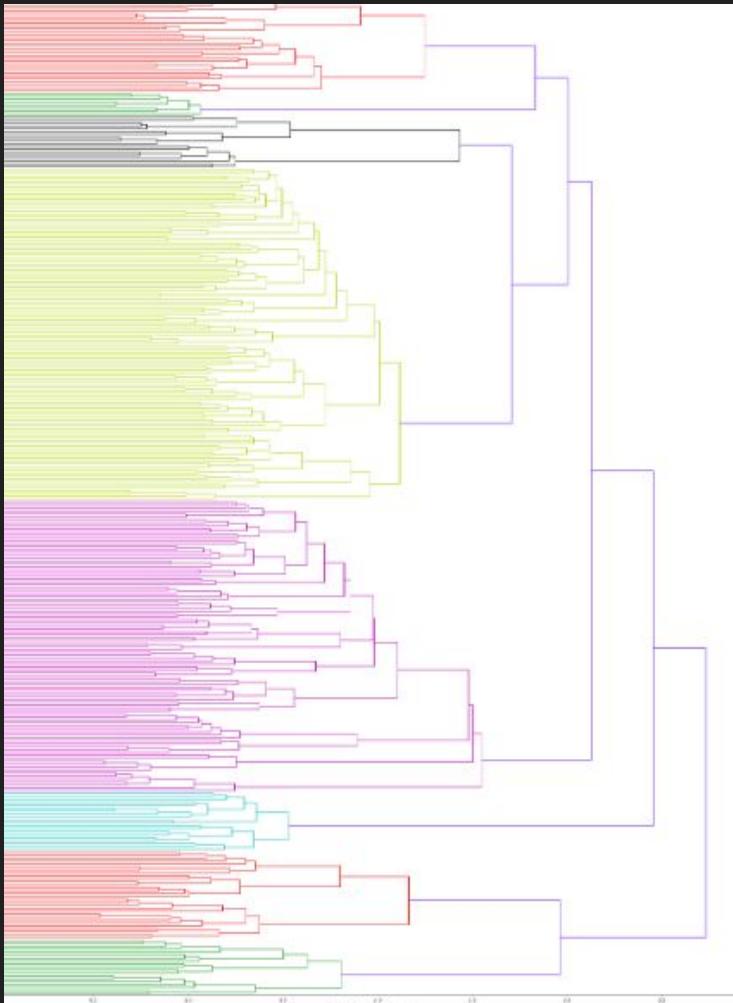


Matrix shape: 432 x 10,865
(432 sections, 10,865 terms)

Matrix shape: 339 x 4,029

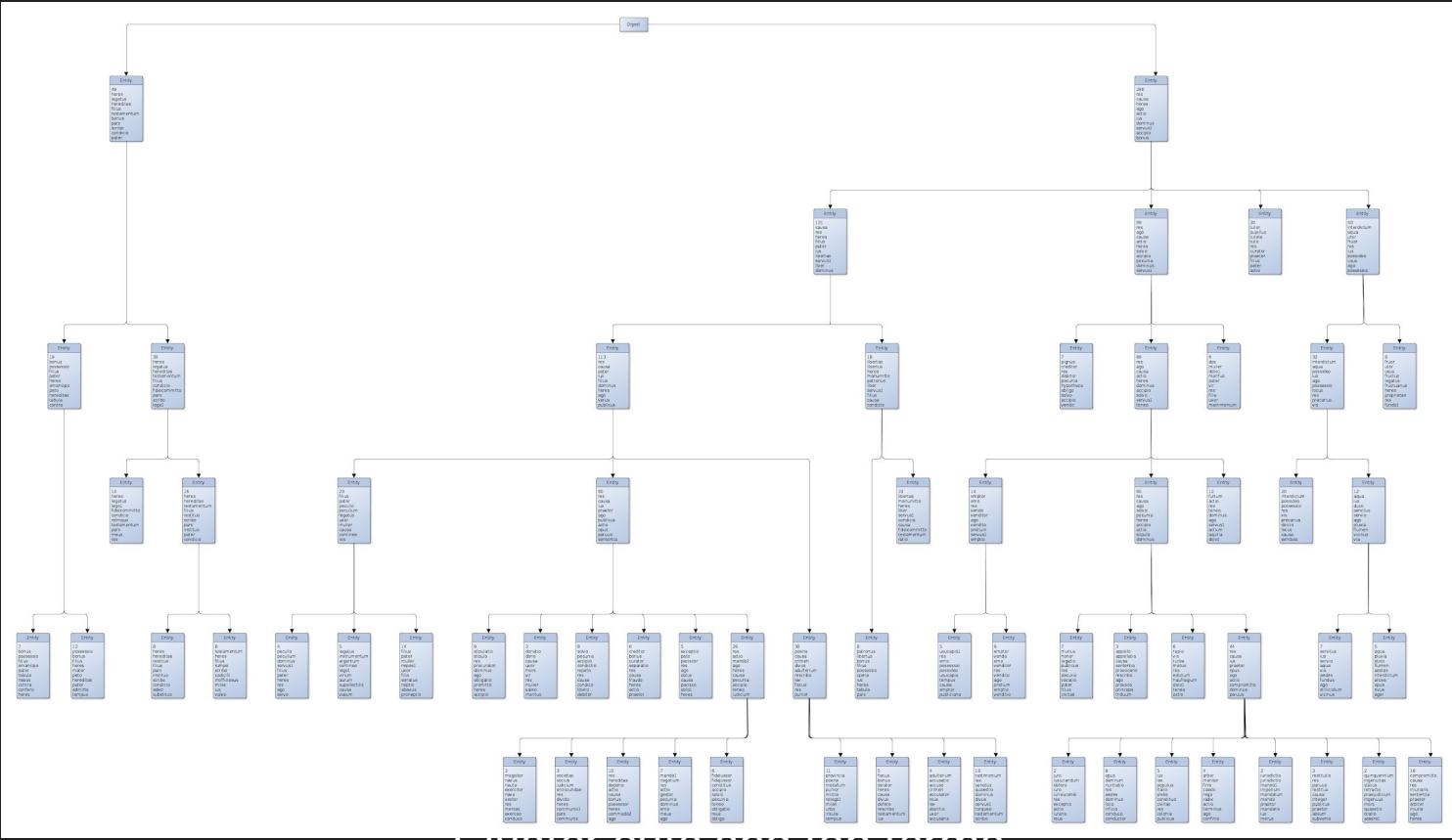


Hierarchical clustering on the tfidf matrix



Threshold (Euclidean distance)	Number of clusters at threshold
3.5	2
3.0	5
2.5	10
2.0	17
1.75	31
1.50	55
1.375	80

https://gitlab.epl.silene.ru/m301/clustered_clustering/NLP_document_clustering



<https://ottawalibrary.ca/branches/otterdown/branches/otterdown/branches/otterdown/documents>

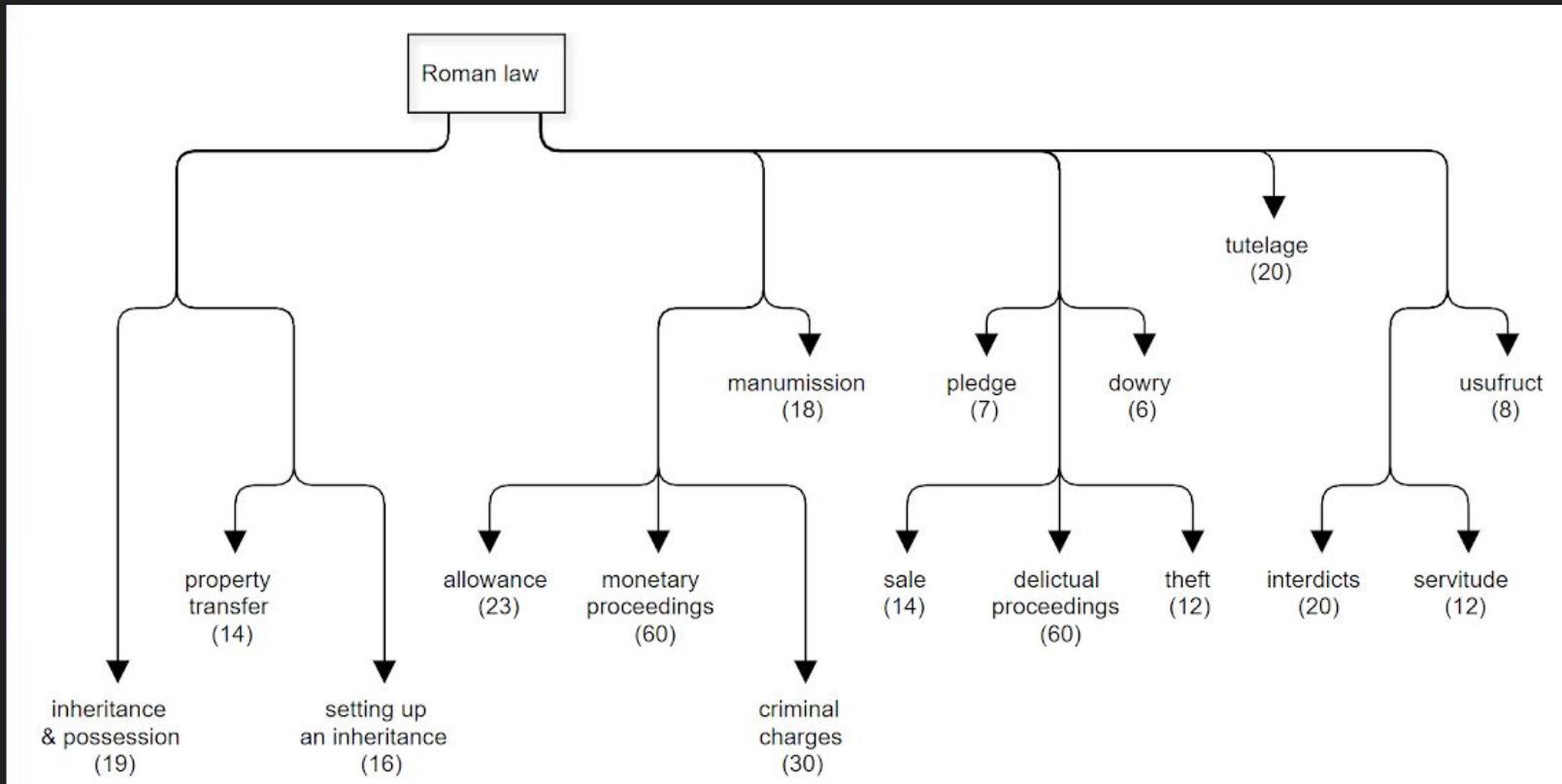
"manumission"

8
libertas
manumitto
heres
liber
servus1
condicio
causa
fideicommitto
testamentum
ratio

"theft"

12
furtum
actio
res
teneo
dominus
ago
servus1
actium
aquia
dolo1

Cluster examples



A simplified empirical conceptual tree-map of Roman law

References

- Burns, M., "Building Better Text Analysis Pipeline for 2018's Old Languages," in Digital classical philology: Ancient Greek and Latin in the digital revolution, edited by Burns, P. J., "Constructing stoplists for historical languages," Digital Classics Online 4:2 (2018): 4-20.
- Burns, B., "Complex lemmatization with the Classical Language Toolkit," presented at the First LiLa Workshop: Linguistic Resources & NLP Tools for Latin on 3 July 2019.
- Klingenberg, G., "Die ROMTEXT-Datenbank," *Informatica e diritto* 4 (1995): 223-232.
- McGillivray, B. and Kilgarriff, A., "Tools for historical corpus research, and a corpus of Latin," on the SketchEngine website. [Accessed on 15 June 2020]
- Mommsen, T. and Krüger, P., eds. *Corpus Iuris Civilis*. Vol 1: *Institutiones. Digesta*. 5th edition. Berlin: Weidmann, 1889.
- Ribary, <https://doi.org/10.5334/johd.17>, "A Relational Database of Roman Law Based on Justinian's Digest." *Informatics* 7 (2020), 44. Available:
- Ribary, M., "A Relational Database of Roman Law Based on Justinian's Digest." *Journal of Open Humanities Data* 6(1), p.5. DOI: <http://doi.org/10.5334/johd.17>
- Ribary, M., "A relational database of Roman law based on Justinian's Digest." figshare. Dataset. <https://doi.org/10.6084/m9.figshare.12333290.v1>
- Ribary, M. pyDigest, A GitLab Repository of Scripts, Files and Documentation. Available online: <https://gitlab.eps.surrey.ac.uk/mr0048/pydigest>
- Riedlberger, P. and Rosenbaum, G., eds. *Amanuensis V5.0*. München: 2020.

Tfidf and dimension reduction

- Term frequency - Inverse document frequency

$$\Rightarrow 432 \times 10,865$$

- Dimension reduction

- Remove short sections (339 remains)

- Select top 50 terms only

$$\Rightarrow 339 \times 4,029$$

TF-IDF

TF-IDF is a measure of originality of a word by comparing the number of times a word appears in a doc with the number of docs the word appears in.

$$\text{TF-IDF} = \text{TF}(t, d) \times \text{IDF}(t)$$

Term frequency
Number of times term t appears in a doc, d

Inverse document frequency
 n ← # of documents

$$\log \frac{1 + n}{1 + \text{df}(d, t)} + 1$$

Document frequency of the term t

Sunoikisis Digital Classics
Session 4. Computational Linguistics



29 October 2020 Leipzig | Monza

YouTube link: <https://youtu.be/zjkyZUpvhAQ>

Computational Linguistics & Deep Learning for Ancient Texts

Thea Sommerschield
DPhil candidate in Ancient History



UNIVERSITY OF
OXFORD

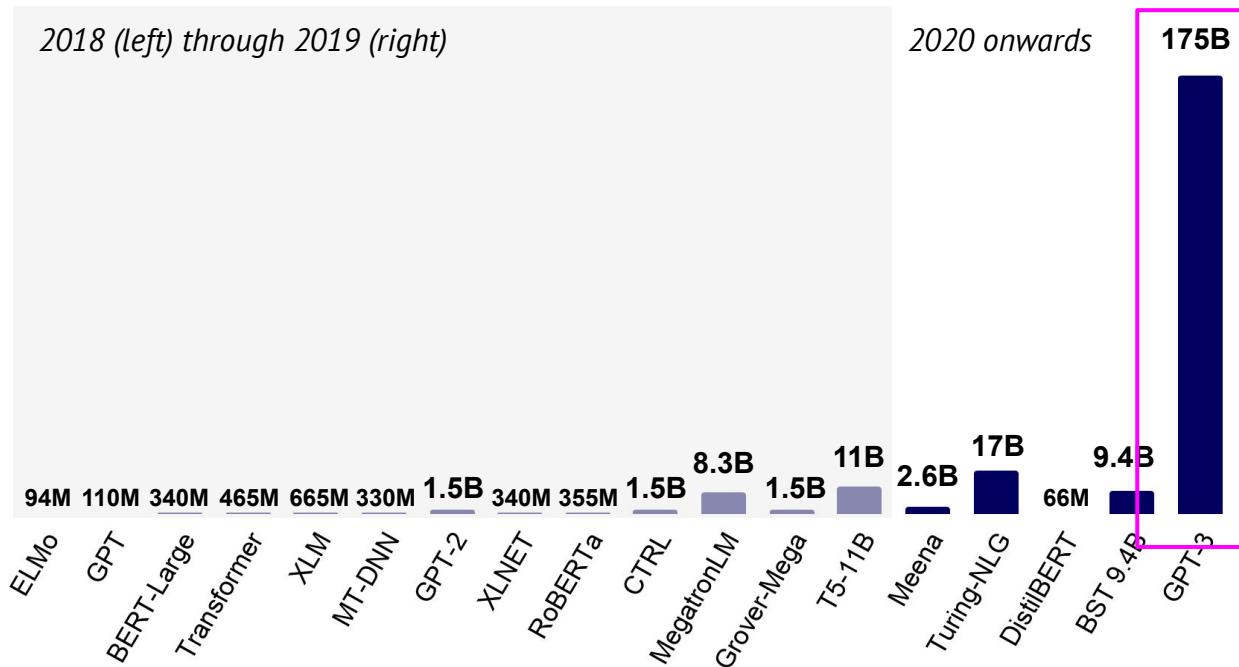
It's a great time to be working on Computational Linguistics and Machine Learning.

- The *2020 State of AI and Machine Learning* report illustrates the current state of artificial intelligence and machine learning, showcasing where the industry is as a whole in 2020 compared to 2019.
- <https://www.stateof.ai/>

Our 2019 Prediction	Grade	Evidence
New natural language processing companies raise \$100M in 12 months.	Yes	Gong.io (\$200M), Chorus.ai (\$45M), Ironscales (\$23M), ComplyAdvantage (\$50M), Rasa (\$26M), HyperScience (\$60M), ASAPP (\$185M), Cresta (\$21M), Eigen (\$37M), K Health (\$48M), Signal (\$25M), and many more!
No autonomous driving company drives >15M miles in 2019.	Yes	Waymo (1.45M miles), Cruise (831k miles), Baidu (108k miles).
Privacy-preserving ML adopted by a F2000 company other than GAFAM (Google, Apple, Facebook, Amazon, Microsoft).	Yes	Machine learning ledger orchestration for drug discovery (MELLODY) research consortium with large pharmaceutical companies and startups including Glaxosmithkline, Merck and Novartis.
Unis build <i>de novo</i> undergrad AI degrees.	Yes	CMU graduates first cohort of AI undergrads, Singapore's SUTD launches undergrad degree in design and AI, NYU launches data science major, Abu Dhabi builds an AI university.
Google has major quantum breakthrough and 5 new startups focused on quantum ML are formed.	Sort of	Google demonstrated quantum supremacy in October 2019! Many new quantum startups were launched in 2019 but only Cambridge Quantum, Rahko, Xanadu.ai, and QCWare are explicitly working on quantum ML.
Governance of AI becomes key issue and one major AI company makes substantial governance model change.	No	Nope, business as usual.

Language models: Welcome to the Billion Parameter club

► Huge models, large companies and massive training costs dominate the hottest area of AI today, NLP.



Source: <https://arxiv.org/pdf/2004.08900.pdf>

Note: The number of parameters indicates how many different coefficients the algorithm optimizes during the training process.

It's a great time to be working on texts from the Ancient World.

- The 2020 *BES Autumn Colloquium* of the *British Epigraphy Society* features 5/10 digital epigraphy initiatives in its yearly Epigraphic Gazetteer.
- <http://www.britishepigraphysociety.org/bes-autumn-colloquium.html>



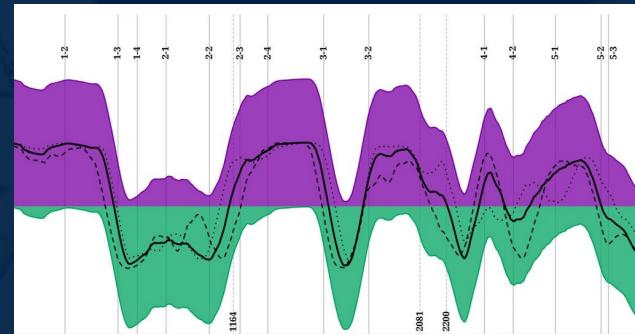


<http://sicily.classics.ox.ac.uk/>

Google Arts & Culture

Bringing AI to ancient languages

<https://artsexperiments.withgoogle.com>



<https://arxiv.org/pdf/1911.05652.pdf>



Machine Translation and Automated
Analysis of Cuneiform Languages

<https://cdli.ucla.edu/> & <https://cdli-gh.github.io/mtaac/>



Neural Decipherment via Minimum-Cost Flow:
from Ugaritic to Linear B

<https://arxiv.org/abs/1906.06718>

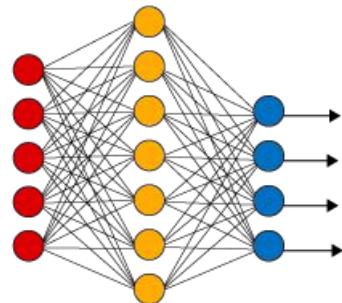


<https://hieroglyphicsinitiative.ubisoft.com/>

But first, some definitions.

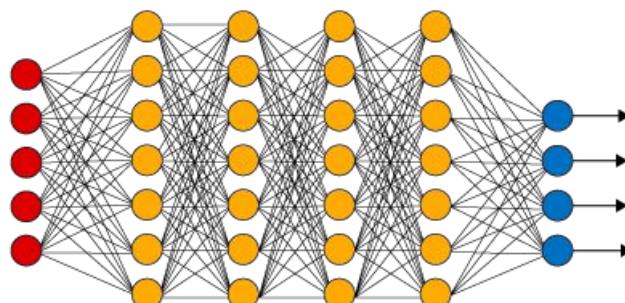
- Machine Learning
- Deep Learning
- Epigraphy (in a bit)

Simple Neural Network



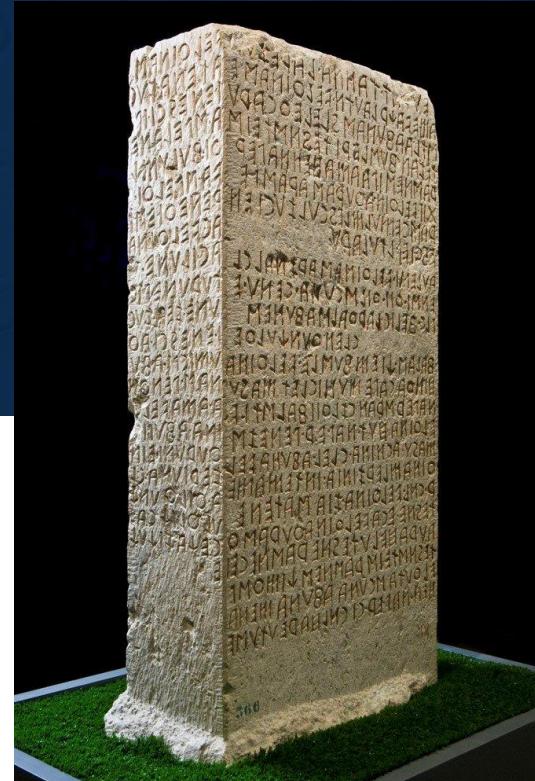
● Input Layer

Deep Learning Neural Network



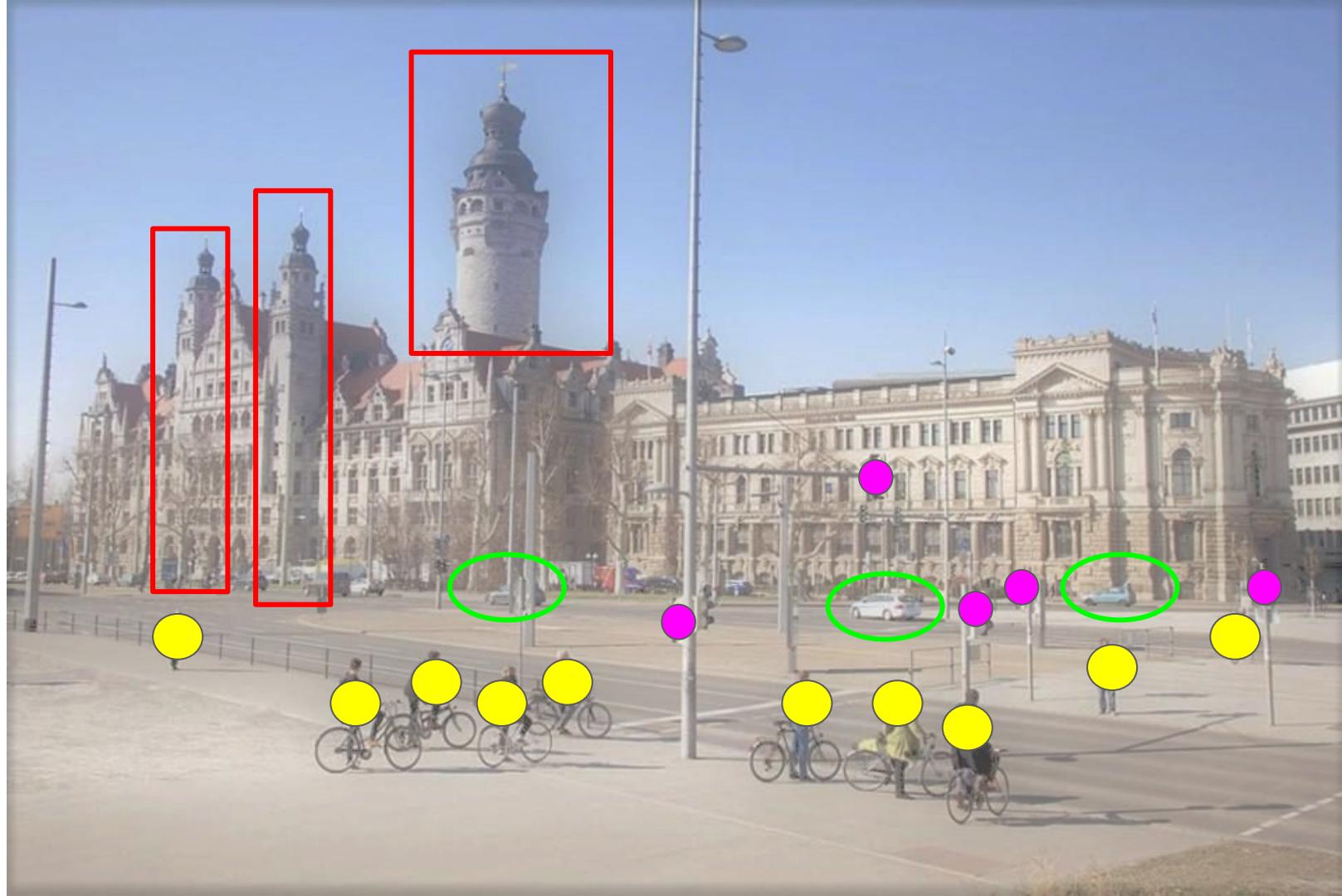
● Hidden Layer

● Output Layer



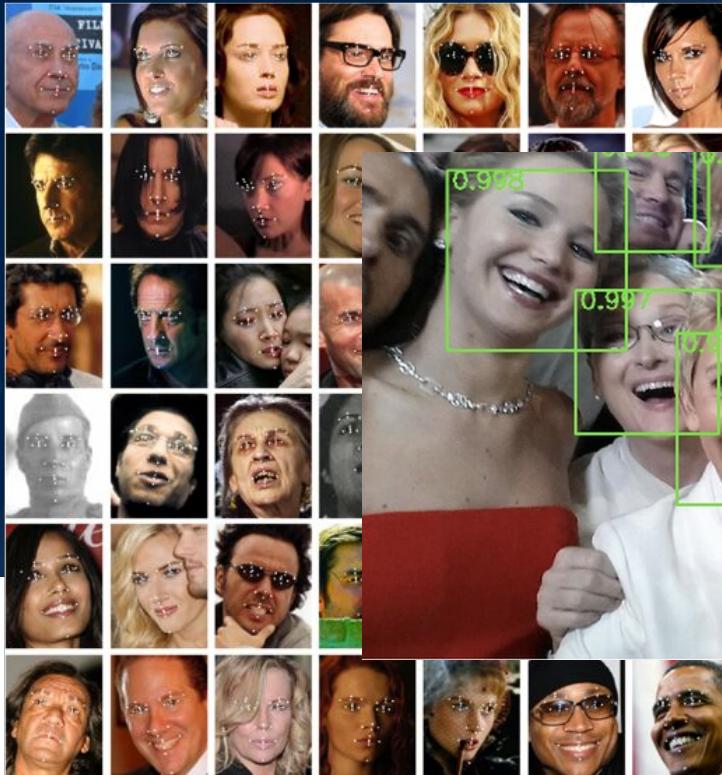
The Cippus Perusinus



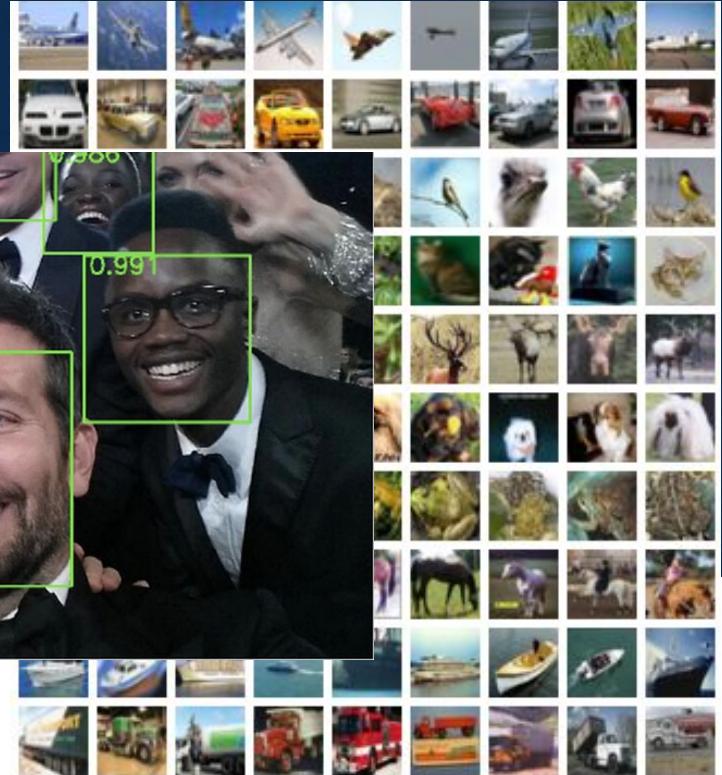


Learning features from data

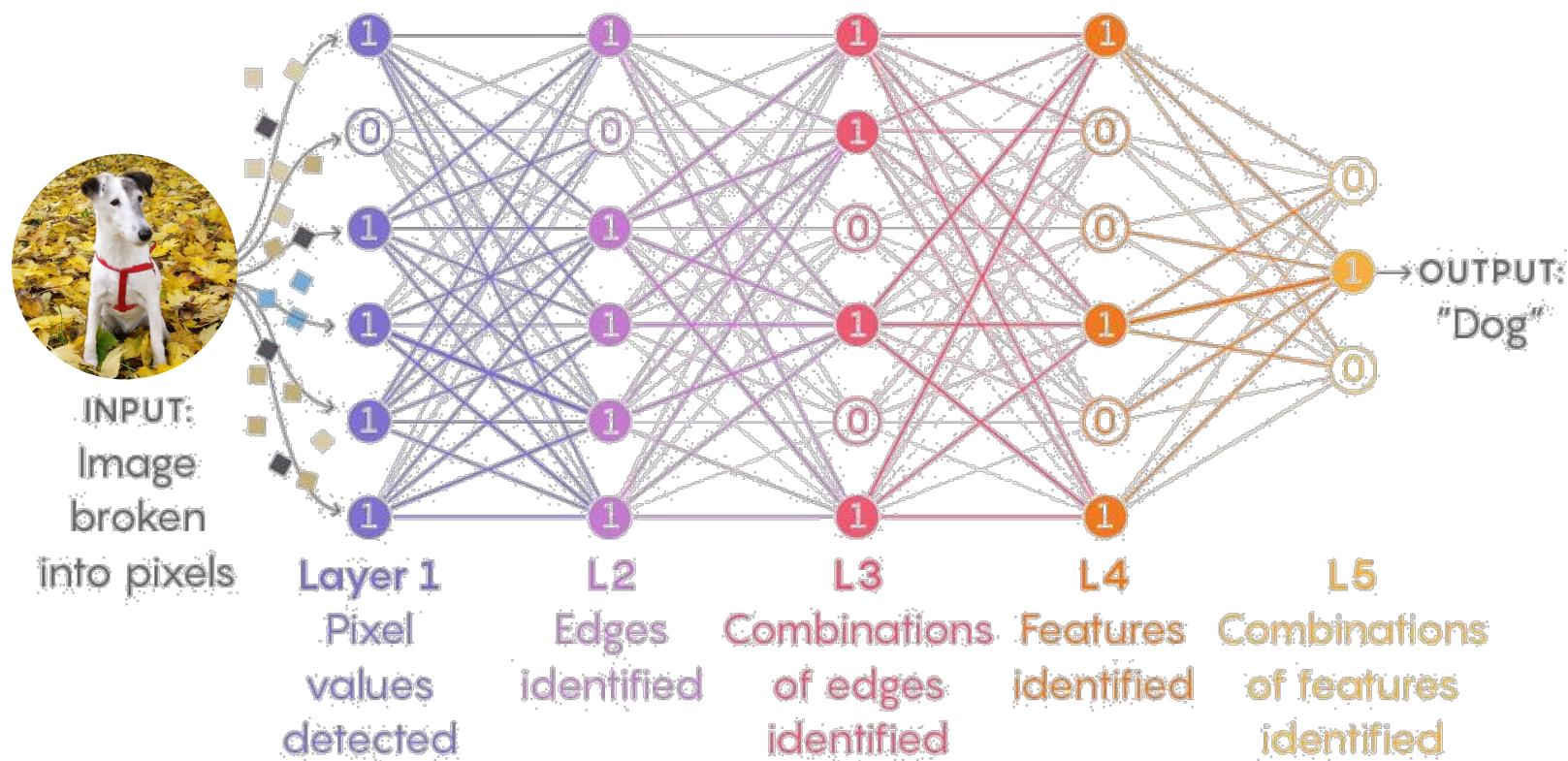
faces



not faces



Neural Networks



Neural Networks

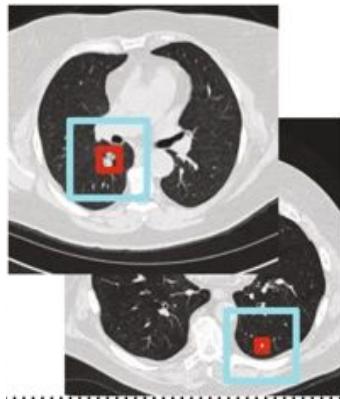
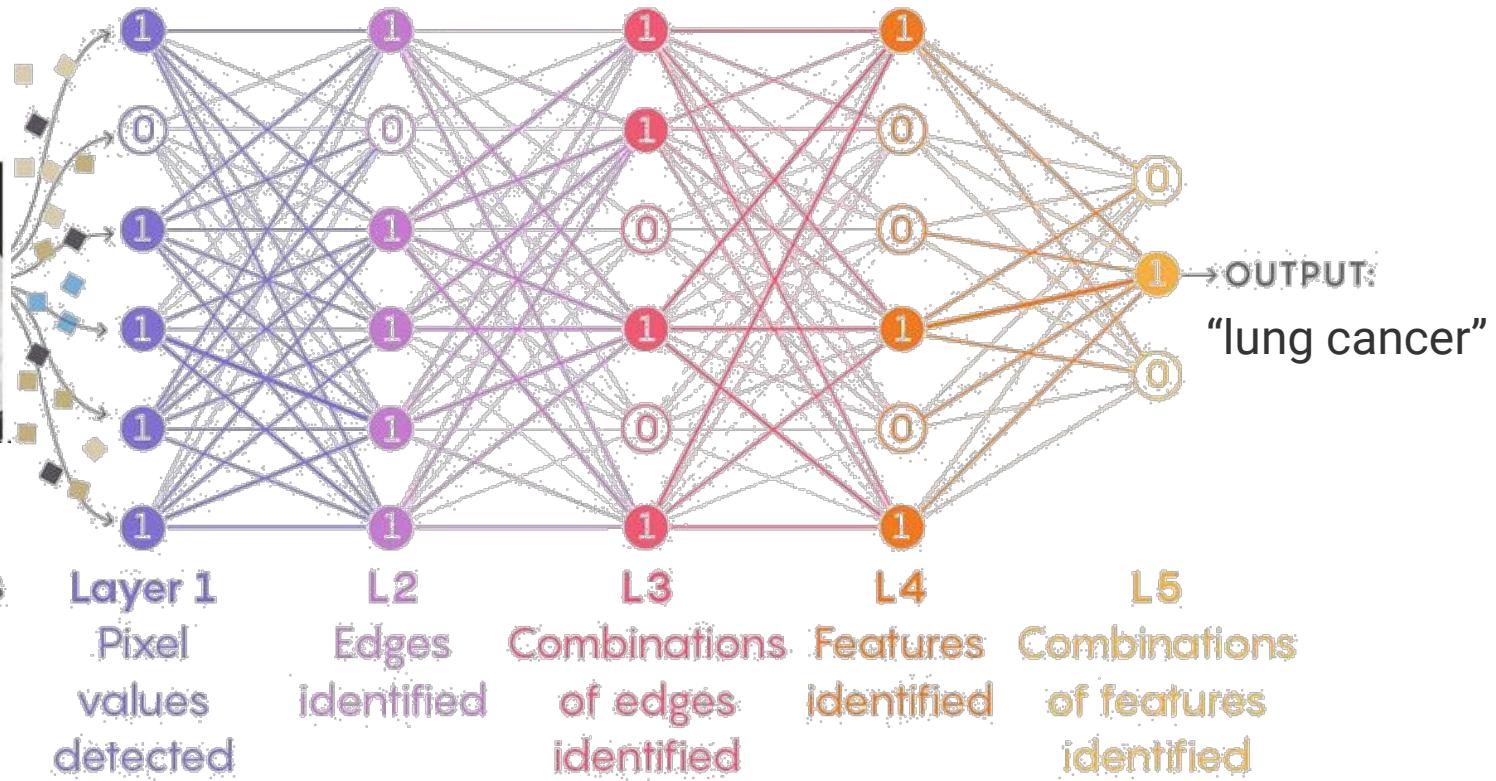


Image
broken
into pixels



“Restoring ancient text using deep learning: a case study on Greek epigraphy”

Yannis Assael*, Thea Sommerschield*, Jonathan Prag



Association for
Computational Linguistics



EMNLP 2019

Conference on Empirical Methods
in Natural Language Processing

Hong Kong, China
November 3–7, 2019

ΤΟΝ ΑΞΙΟΛΟΓΩΣ ΤΑ ΤΟΝ ΑΜΟΙ
ΚΤΥΟΝ ΑΓΕΝΟ ΜΕΝ ΟΝΤΥ ΘΙ
ΑΔΙ ΜΑΡΟΥΑ ΠΙΟΝ ΑΚΗ
ΤΙΟΝ ΛΟΥΚΙΟΝ ΝΕΙΚΟΠΟ
ΧΕΙΓΗΝ ΔΑΡΤΕ ΗΘΟΥ ΣΕΓΗΝ
< ΕΙΑΝ ΚΑΙ ΔΑΤΗΝ ΣΕ ΠΟΥ ΔΗΝ
ΝΕΝΕΔΕΙ ΕΙΣ ΤΟΥ ΠΕΡΤΗΣ ΣΕ
ΜΝΟΠΕΠΟΣ ΣΤΟ ΛΓΟΝ ΝΟΣΤΩΝ
ΜΕΓΑΛΩΝΤΥ ΘΙΘΟ ΛΠΡΟΙΚΑ
ΠΡΕΣΒΕΥΣ ΣΑΝ ΣΑΟΙ ΚΕΙΟΣ
ΤΕΡΦΤΗΝ ΣΕ ΣΑΜΕΝΟΝ
ΛΑΤΟΥΑΝΔΑ ΧΟΣ Λ Α Ν Α
ΣΤΑΣ

Aims:

- A fully automated aid to the epigraphist's restoration task: a ML model.
- An epigraphic dataset of machine actionable text to train the model.
- Design a model which can be used by all disciplines dealing with ancient texts and which applies to any language.



UNIVERSITY OF
OXFORD

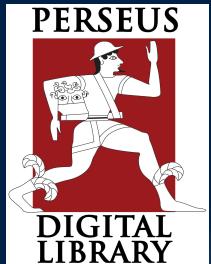
Why ancient Greek epigraphy?

1. The variability of its contents and context within this dataset makes it an excellent challenge for Machine Learning.
2. Increased availability of digitised data.
3. Personal interest/expertise.



UNIVERSITY OF
OXFORD

Where can I find thousands of digitised Greek texts?



THE PACKARD HUMANITIES INSTITUTE
LOS ALTOS, SANTA CLARITA, CAMBRIDGE



A system is as good only as the data it collects!

1. We need a lot of data.
2. We need a lot of clean data.



UNIVERSITY OF
OXFORD

Regions : Attica (IG I-III)

Agora XVI 2[1] ← Agora XVI 1[5] Agora XVI 2[2] →

[1]

Att. — Athens: Agora — stoich. — c. 450 a. — IG I³ 16 — Hesp. 14.1945.82,2; 15.1946.246,77 — SEG 10.11

frq. a

{²reliqiae nonnullae}²

1 [.] ταὶ ἐὰμ[— — — — — — —]
[.] σθο[.] ατ[— — — — — —]
[...] Λ! Ο#?
[.....] τα καὶ [— — — — — —]
5 [.. 7 ..] εμ[— — — — — —]
[.. 6 ..] ΟΛΕ[
[..] τὲν βολὲ[ν — — — — —]
[φ]ρόραρχον [— — — — φ]
[ρ]οφεν πεντέ[κοντα — — — —]
9 οἱ ἐν ἀκροπόλι[ει — — — — —]
φρόραρχον α[— — — — — —]
[.] τὸν δὲ φρ[όραρχον — — — τ]-
ριάκοντα ἔ[τε γεγονότας(?) — — βο]-
εθείαι [.] ακ[— — — — — —]

Machine readable text:

1. **Text cleanup:** compute character frequencies, standardise alphabet, strip human annotations and inconsistencies.
2. **Text processing:** ‘-’ for missing characters and ‘?’ for characters to predict; match the number of ‘-’ with those conjectured by epigraphists.



UNIVERSITY OF
OXFORD

[— — — Ἄ]πολλωδώρου
 [Εὐβ]οῖδος
 10 — — ὥνιος Ἀττίνου
 [Εύμ]ενείας
 — — — Διονοσίου {²⁶Διονυσίου}²⁶
 [Εύμ]ενείᾳ[ς]
 — — — — —
 15
 c — — —
 I — — — — εστράτου
 — — — εύς
 [— — — — —]!τουτων
 [(τῶν) ἔξ Ἀββου κ]ώμης
 5 — — — Ἀπολλωνίου
 — — — νος
 — — — — ου Μυσ[ός]

ἀπολλωδώρου εὐβοϊδος --ώνιος ἀττίνου
 εὐμενείας --- διονοσίου εὐμενείας
 -----εστράτου ---εύς -----ιτουτων
 τῶν ἔξ ἀββου κώμης -- **ἀπολλωνίου**
 ---νος -----ου μυσός

?????????ou

The PHI-ML training set:

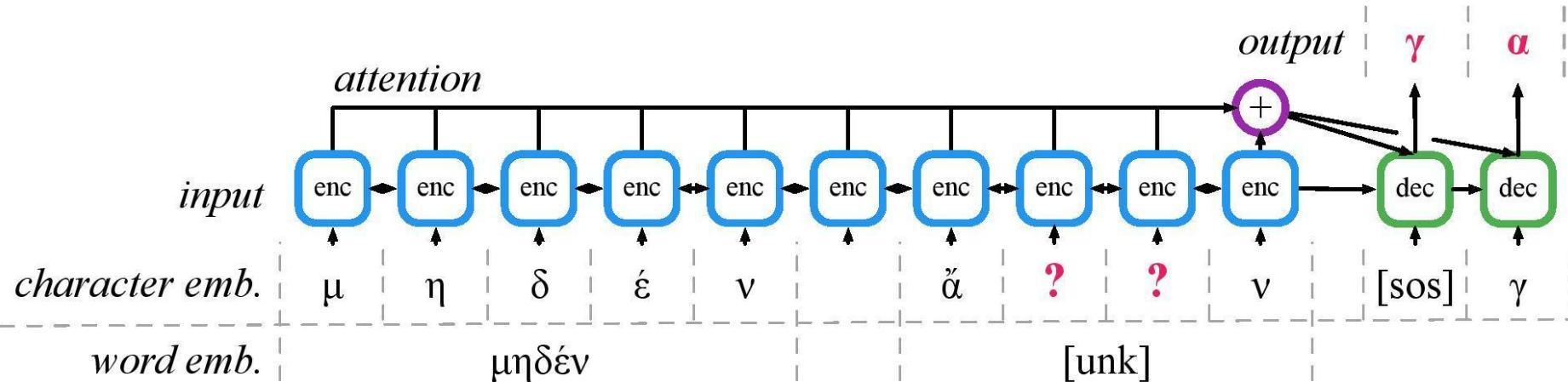
Split	Inscriptions	Words	Chars
Train	34,952	2,792k	16,300k
Valid	2,826	211k	1,230k
Test	2,949	223k	1,298k

Statistics for the PHI-ML corpus.

Architecture design:

- **Input:** sequence of characters
- **Output:** restoration predictions
- **seq2seq** architecture turning input sequence into output sequence
- **LSTM** recurrent neural network
- Works at a **word and a character level**

Pythia:



Pythia BI-WORD processing the phrase **μηδέν** **ᾰγαν.**
 The letters “γα” are the characters to be predicted, and are annotated with ‘?’.
 Since ᾰ??v is not a complete word, its embedding is treated as unknown ('unk').
 The decoder outputs correctly “γα”.

Ablation study:

Method	CER	Top-20
Ancient Historian	57.3%	—
LM Philology	68.1%	26.0%
LM Philology & Epigraphy	65.0%	28.8%
LM Epigraphy	52.7%	47.0%
PYTHIA-UNI	42.2%	60.6%
PYTHIA-BI	32.5%	71.1%
PYTHIA-BI-WORD	30.1%	73.5%

Predictive performance on PHI-ML.

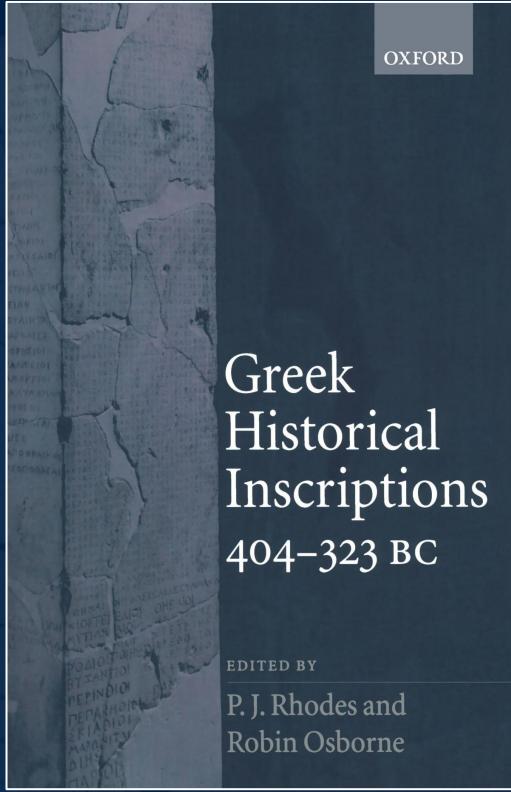
Text restoration demo:

[https://colab.research.google.com/drive/16RfCpZ
Lm0M6bf3eGIA7VUPclFdW8P8pZ](https://colab.research.google.com/drive/16RfCpZLm0M6bf3eGIA7VUPclFdW8P8pZ)



UNIVERSITY OF
OXFORD

IG II² 116



UNIVERSITY OF
OXFORD

θεοί ἐπὶ νικοφήμῳ ἄρχοντος συμμαχίᾳ ἀθηναίων καὶ θετταλῶν εἰς τὸν ἀεὶ χρόνον. ἔδοξεν τῇ βουλῇ καὶ τῷ δῆμῳ λεωντὶς ἐπρυτάνευεν χαιρίων χαριναύτῳ φαληρεὺς ἐγραμμάτευεν ἄρχιππος ἀμφιτροπῆθεν ἐπεστάτει δωδεκάτει τῆς πρυτανείας ἐληκεστίδης εἶπεν περὶ ὧν λέγει οὐσιν οἱ πρέσβεις τῶν θετταλῶν ἐψηφίσθαι τῷ δῆμῳ δέχεσθαι τὴν συμμαχίαν τύχῃ ἀγαθῇ καθὰ ἐπανγέλλονται οἱ θετταλοί. εἶναι δὲ αὐτοῖς τὴν συμμαχίαν πρὸς ἀθηναίος εἰς τὸν αἰεὶ χρόνον. εἴναι δὲ καὶ τοὺς ἀθηναίων συμμάχους ἀπαντας θετταλῶν συμμάχος καὶ τὸς μετταλῶν ἀθηναίων. ὅμοσαι δὲ ἀθηναίων μὲν τὸς στρατηγὸς καὶ τὴν βολὴν καὶ τὸς ἵππαρχος καὶ τὸς ἵππαρχος τόνδε τὸν ὄρκον βοηθήσω παντὶ σθένει κατὰ τὸ δυνατόν ἐάν τις ἦτι ἐπὶ τὸ κοινὸν τὸ θετταλῶν ἐπὶ πολέμῳ ἢ τὸν ἄρχοντα καταλύει ὃν εἴλοντο θετταλοί ἢ μύραννον καθὶ στῇ ἐν θετταλίαιεπομύναι δὲ τὸν κόμιμον ὄρκον. ὅπως δὲ καὶ θετταλοί ὅμόσωσι τῇ πόλει ἐλέσθαι τὸν δῆμον πέντε ἄνδρας ἐν ἀθηναίων ἀπάντων οἵτινες ἀφικόμενοι εἰς θετταλίαν ἐξορκώνοσιν ἀγέλαον τὸν ἄρχοντα καὶ τὸς πολεμάρχος καὶ τὸς ἵππαρχος καὶ τὸς ἵππαρχος καὶ τὸς ἱερομνήμονας καὶ τοὺς ἄλλους ἄρχοντας ὁπόσοι ὑπὲρ τὸ κοινὸν τὸ θετταλῶν ἄρχοσαν τόνδε τὸν ὄρκον βουθέτω παντὶ σθένει κατὰ τὸ δυνατόν ἐάν τις ἦτι ἐπὶ τὴν πόλιν τὴν ἀθηναίων ἐπὶ πολέμῳ ἢ τὸν δῆμον καταλύει τὸν ἀθηναίων ὅμοσαι δὲ καὶ τὸς πρέσβεις τὸς τῶν θετταλῶν ἐν τῇ βολῇ τὸς ἐπιδημοντας ἀθήνησιν τὸν αὐτὸν ὅλον τὸν δὲ πόλεμον τὸν πρὸς ἀλέξανδροντὸν μὴ πόλεμαν καταθύσασθαι τοῖς θετταλοῖς ἀνευ ἀθηναίων τοῖς ἀθηναίοις ἄρχοντο ἄρχοντος καὶ τοῦ κοινοῦ τῶν θετταλῶν. ἐπαινέσαι δὲ ἀγέλαον τὸν ἄρχοντα τὸν στρατηγὸν τῶν θετταλῶν ὅτι εὗ καὶ προθύμως ἐπιμελεσασθαι περὶ ὧν αὐτοῖς ἢ πόλις ἐπηγγείλατο ἐπαινέσαι δὲ καὶ τὸς πρέσβεις τῶν θετταλῶν τὸ ἄρχοντας καὶ καλέσαι αὐτὸς ἐπὶ ξένια εἰς τὸ πρυτανεῖον εἰς αὔριον. τὴν δὲ στήλην τὴν πρὸς ἀλέξανδρον ἀνθελλων τὸς ταμίας τῆς θεοτητὸς τὰς συμμαχίας. τοῖς δὲ πρέσβεις δοναι τὸν ταμίαν τοῦ δῆμο εἰς ἐφόδια δδ δραχμὰς ἐκάστωι. τὴν δὲ συμμαχίαν τίδε διαγράψαι τὸν γραμματέα τῆς βολῆς ἐνστήλῃ λιθίνῃ καὶ στῆσαι ἐν ἀκροπόλει εἰς δὲ τὴν ἀναγραφὴν τῆς στήλης δοναι τὸν ταμίαν το δῆμο 0 δραχμὰς εἶναι δὲ καὶ θειότητον τὸν ἐρχιέα ὡς λέγοντα ἄριστα καὶ πράττοντα ὅ τι ἀν δύνηται ἀγαθὸν τῶιδήμῳ τῷ ἀθηναίων καὶ θετταλεῖς ἐν τῷ τεταγμένῳ.

Sample restoration of the inscription *IG II² 116*.
Restorations are in blue when correct, purple when incorrect.

Results:

1. Our experimental results illuminate the ways Pythia can assist, guide and advance the epigraphist's task.
2. The combination of Machine Learning and Epigraphy can impact meaningfully the study of inscribed textual cultures, ancient and modern.



UNIVERSITY OF
OXFORD

Impact and significance:

1. Track textual connections and correlations with computational methods.
2. Design educational tools to explore the data.
3. Make standardised and organised data more accessible.

For more information:

- Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP):
<https://www.aclweb.org/anthology/D19-1668/>
- DeepMind blog post:
<https://deepmind.com/research/publications/Restoring-ancient-text-using-deep-learning-a-case-study-on-Greek-epigraphy>
- University of Oxford blog post:
<http://www.ox.ac.uk/news/arts-blog/restoring-ancient-greek-inscriptions-using-ai-deep-learning>
- Financial Times article:
https://www.ft.com/content/2b72ed2c-907b-11ea-bc44-dbf6756c871a?FTCamp=engage%2FCAPI%2F%2FChannel_signal%2F%2FB2B





Thank you!

Thea Sommerschield

DPhil candidate in Ancient History

thea.sommerschield@classics.ox.ac.uk



UNIVERSITY OF
OXFORD



A natural language processing ‘pipeline’ of the papyri

Alek Keersmaekers

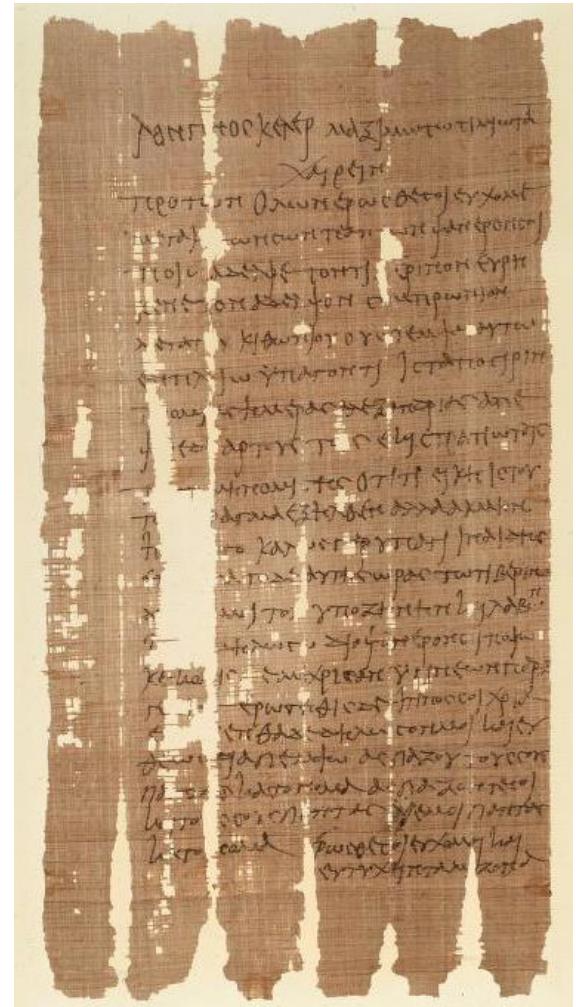


University of Leuven
RU Quantitative Lexicology and Variational Linguistics



Introduction: the papyri

- “Documentary” texts from Greek/Roman Egypt (3rd century BC – 8th century AD)
- Letters, petitions etc.: closer to everyday language
- More sociolinguistic variation: gender, social class, age, mother tongue etc.
- 4.5 million tokens
- Current linguistic work on papyri: mostly small-scale (e.g. Porter & O'Donnell 2010: 3341 words)



Towards an automatic analysis ‘pipeline’

- All (published) texts are available in a computer-readable (XML) format, but lacking any sort of tokenization, linguistic annotation etc.

```
▼<div xml:lang="grc" type="edition" xml:space="preserve">
  ▼<div n="r,ctr" type="textpart">
    ▼<ab>
      <lb n="1"/>
      ὁ ἐπιφ
      <supplied reason="lost">a</supplied>
      <unclear>v</unclear>
      ἐστατος Φωνην βασιλεύς Βλεμμύων
```

- Supplying missing information via an NLP ‘pipeline’



Main problems

- Greek: highly inflectional, free word order, large rate of elliptic constructions
- Incompletely preserved texts (i.e. ‘gaps’ in some sentences)
- Not a homogeneous spelling
- Mismatch between training corpus and test data



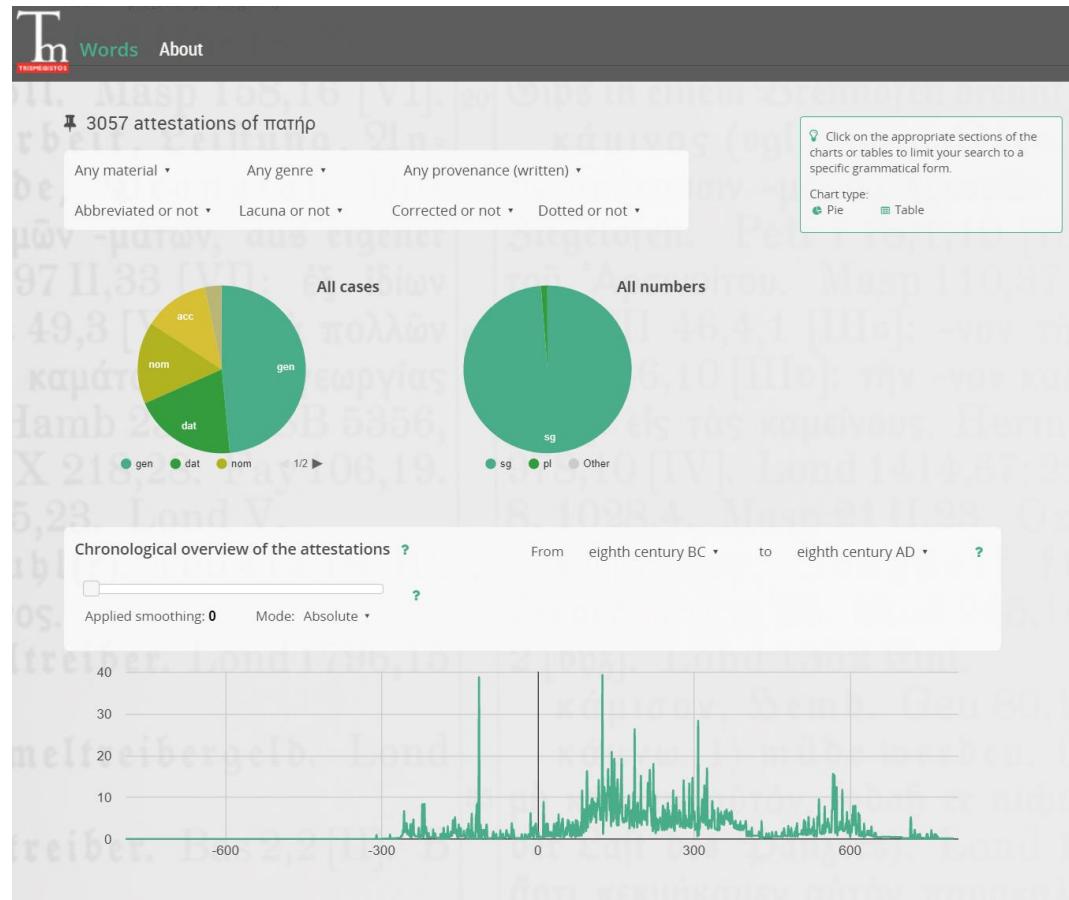
The Trismegistos project

- Collaboration with Trismegistos (Ancient History, KU Leuven) through Mark Depauw (co-supervisor of PhD project)
 - 800BC-800AD
 - ‘Metadata’ (e.g. date, provenance, onomastic data)...
 - ... of all non-literary texts (including papyri)
- Linked open data
- *Trismegistos Words*: linguistic database
- Mutually beneficial for linguistics and ancient history



Trismegistos Words

- <https://www.trismegistos.org/words>
- See also <https://github.com/alekkeersmaekers/duke-nlp>



Tokenization

- Tokenization: rule-based approach (spaces and punctuation), some persistent problems (capitalization, multiple separate words written together) resolved semi-automatically
- Representing different text versions (original/editorial) under same token
- Possible to link to existing info in Trismegistos databases (People, Places, Text Irregularities)

ὅπως καὶ γὰρ σοὶ(∗) δέ, ἔχω παῖδαν(∗) καὶ κάγώ(∗).

^ r,ctr.5.1. σὺ
^ r,ctr.5.1. παῖδα
^ r,ctr.5.1. καὶ ἔγώ

```
<sentence id="5">
    <word id="1" row="5" wordNum="5" hand="m1" regularized="ὅπως" original="ὅπως" comple
    <word id="2" row="5" wordNum="6" hand="m1" regularized="καὶ" original="καὶ" comple
    <word id="3" row="5" wordNum="7" hand="m1" regularized="γὰρ" original="γὰρ" comple
    <word id="4" row="5" wordNum="8" hand="m1" regularized="σὺ" original="σοὶ*" comple
    <word id="5" row="5" wordNum="9" hand="m1" regularized="δέ" original="δέ" completen
    <word id="6" row="5" wordNum="10" hand="m1" regularized="," original="," completen
    <word id="7" row="5" wordNum="11" hand="m1" regularized="ἔχω" original="ἔχω" comple
    <word id="8" row="5" wordNum="12" hand="m1" regularized="παῖδα" original="παῖδαν*"
    <word id="9" row="5" wordNum="13" hand="m1" regularized="καὶ" original="καὶ" comple
    <word id="10" row="5" wordNum="14" hand="m1" regularized="κάγω" original="κάγώ*" com
    <word id="11" row="5" wordNum="15" hand="m1" regularized="ἔγώ" original="|extra|*"
    <word id="12" row="5" wordNum="16" hand="m1" regularized="." original="." completen
</sentence>
```

Morphological tagging: introduction

- For Ancient Greek: part-of-speech (e.g. noun, verb) + morphology (e.g. genitive, active, ... > 33 possible features)
- Typical approach: using lexical probabilities from corpus + N preceding words (e.g. *I can do this*)

Morphological tagging: problems

- Number of possible forms for a given lemma is much higher than for English (e.g. λελυκώς: participle, singular, perfect, active, masculine, nominative)
 - Data sparsity problem: typically inflectional languages are analyzed first formally (in this case via *Morpheus*)
 - Special machine learning algorithms are needed to deal with morphologically complex languages
- Syntax/morphology highly interrelated (e.g. nominative/accusative)
- Irregular spelling > editorial corrections
- Unknown words (e.g. names): extension Morpheus, information Trismegistos databases

τύχανω
(Show lexicon entry in LSJ Middle Liddell Slater Autenrieth) ([search](#))

τύχῃ verb 3rd sg aor subj act

τύχῃ verb 2nd sg aor subj mp

[Word frequency statistics](#)

Morphological tagging: results

- Overall accuracy: about 94.7% (*RFTagger, Mate*)
- Some morphological attributes are easier to tag than others > less ambiguous forms

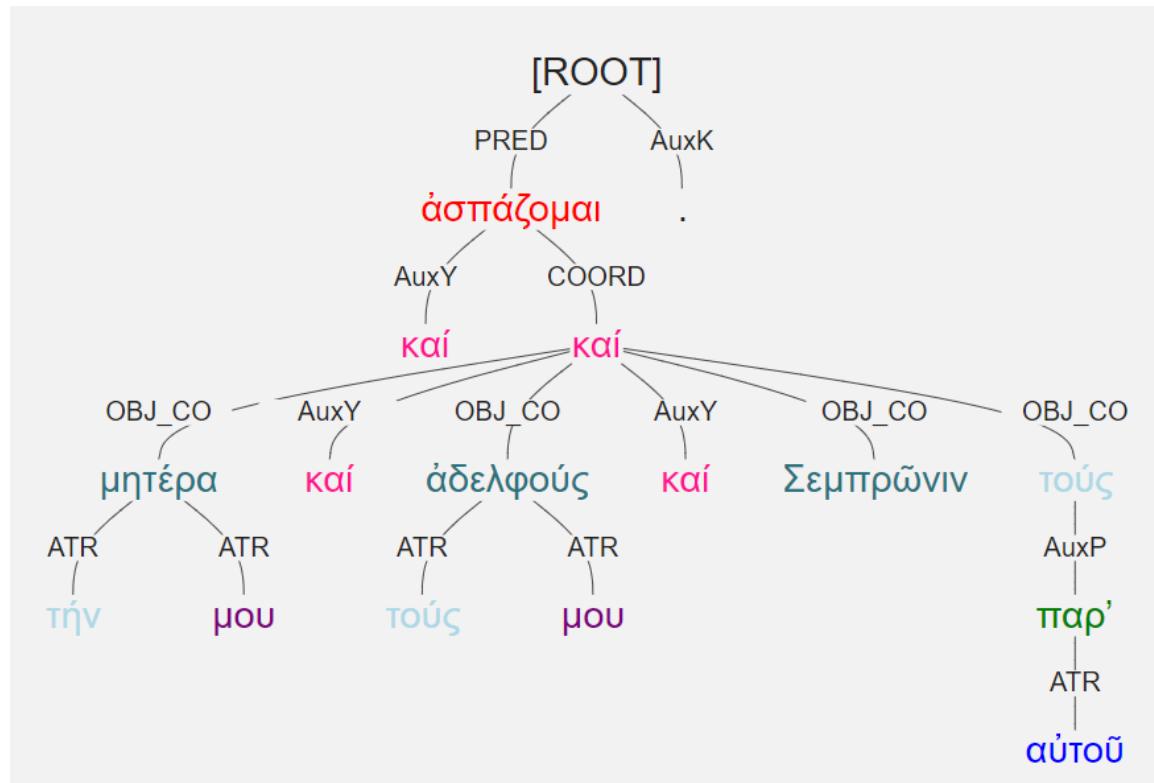
	Median accuracy
Derivative category	0.996
Part-of-speech	0.994
Number	0.990
Voice	0.989
Tense	0.985
Degree	0.974
Case	0.974
Person	0.963
Mood	0.959
Gender	0.951

Morphological tagging: remaining problems

- How to resolve the ~5% remaining problems?
- Some obvious steps in NLP are moving to ‘better’ algorithms (e.g. deep learning) and expanding the training data
- However, taking a deeper look at remaining problems: almost all of them syntactical (e.g. long distance between verb and argument: **ἔστι** γὰρ τὸ πλῆθος τοῦ ἀργυρίου οὐκ **όλιον** “Since the sum of the money **is** not **small**”)
- Smaller category of problems: related to world knowledge (e.g. past verbs on –ov: ambiguous between 1 singular and 3 plural)

Syntactic parsing

- For Ancient Greek: dependency parsing
- Two predictions: syntactic head + label



Graph-based parsing

- Initially, all words have a possible dependency link with each other words
- Machine learning is used to determine which dependency link is the most plausible

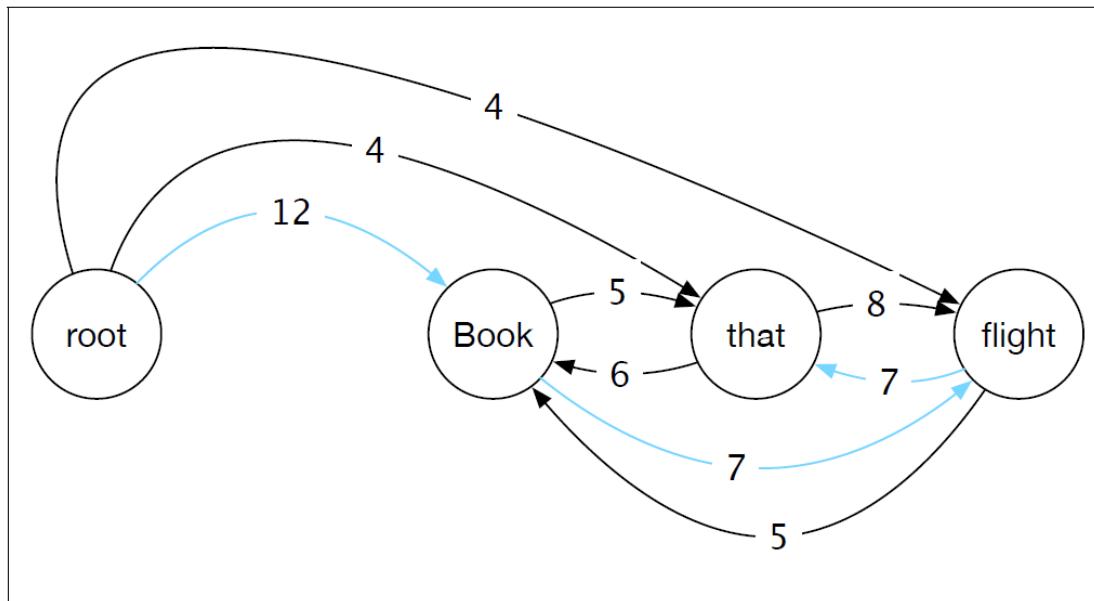


Figure 13.12 Initial rooted, directed graph for *Book that flight*.

Syntactic parsing: main problems (1)

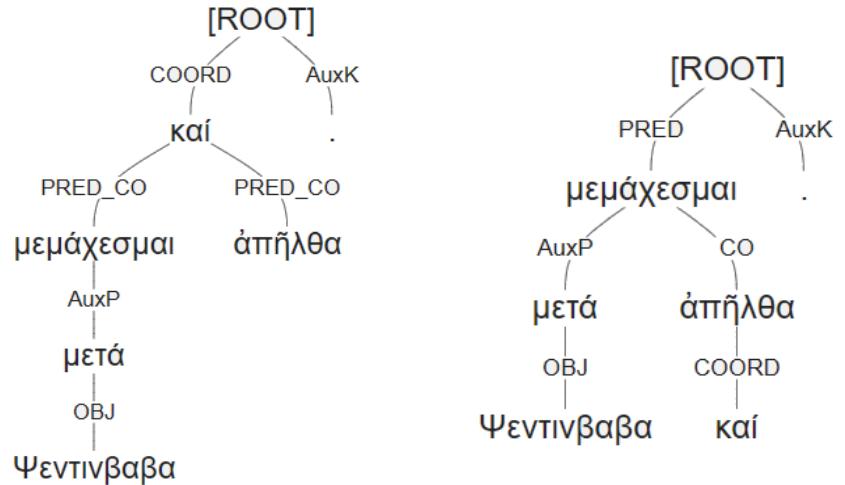
- Major Greek treebanks (e.g. AGDT, Gorman, PROIEL): high rate of inconsistencies (PROIEL even completely different style)

➤ Homogenization

- Manual
- Rule-based
- In the future: machine learning techniques

- More difficult structures: e.g. coordination

➤ Modifying annotation format for coordination structures



Syntactic parsing: main problems (2)

- Damaged text

σε(_). ίμα[τιο]ν χρωμ[ά]τ[ιον ε]ύσχημον. πληρόθητι(_). δὲ
5 [.....]α τῆς οἰκία[ς ...] οσθης τῇ τιμῇ τοῦ ίμα-
[τίο]ν. [..]ω. ε. σ. η[....]ν τὸν Αρσινοείτην. μ-



- “Partial parsing” (i.e. feeding such a sentence to the parser like a regular sentence) > some dependencies are still possible to make (e.g. articles with nouns etc.), with OK parsing quality

Syntactic parsing: results

- Stanford’s Graph-Based Neural Dependency Parser: accuracy rate of 84.5% (head+label correctly predicted)
- However, not all ‘mistakes’ are actually wrong
- ‘Real’ mistakes: often related to semantics and world knowledge

Error type	Frequency
Grammatical mistake	277/500 (55%)
Consistency issue	130/500 (26%)
Annotation error in the test data	48/500 (10%)
Technical (multiple nodes without head)	20/500 (4%)
Damaged text	13/500 (3%)
Ambiguous sentence structure	12/500 (2%)

Syntactic parsing: results by genre

Genre	N tokens	Mean LAS	Median LAS	Std Dev
Papyri	17,609	0.845	-	-
Religion	8,166	0.881	0.873	0.010
Example sentences	1,637	0.870	-	-
Biography	1,445	0.832	0.832	0.001
Epistolography	255	0.828	0.803	0.037
History	14,393	0.825	0.825	0.029
Oratory	4,482	0.822	0.818	0.025
Narrative	2,019	0.804	0.820	0.066
Dialogue	2,329	0.798	0.782	0.025
Philosophical Dialogue	2,431	0.790	0.790	0.040
Philosophy & Science	1,608	0.751	0.758	0.024
Poetry	622	0.740	0.726	0.058

Semantic role labeling

- 29 semantic relations (Pedalion grammar: en.pedalion.org)
- E.g. ἦλθε εἰς τήλιν > “direction”
- Features: form (e.g. εἰς+accusative), lemma, animacy, part-of-speech, morphology, syntactic relation, **word vectors/embeddings** (representation of the ‘meaning’ of a word in a series of numbers)
- Training examples: about 12500, from various treebanks (mostly Pedalion treebanks)

Semantic role labeling: results

- Overall accuracy: about 75.7% for all Greek texts

	Accuracy
Religion	0.838 (932/1112)
Documentary	0.809 (1332/1646)
History	0.765 (1439/1881)
Drama	0.751 (1091/1453)
Narrative	0.751 (2019/2689)
Rhetorical	0.723 (1086/1503)
Philosophy	0.714 (1076/1506)
Epic and lyric poetry	0.687 (485/706)

	F1
recipient (1289)	0.909
material (22)	0.842
direction (1006)	0.840
intermediary (16)	0.815
source (803)	0.797
agent (364)	0.785
time (943)	0.777
manner (1596)	0.775
degree (295)	0.768
companion (424)	0.765
location (1436)	0.752
possessor (127)	0.739
duration (221)	0.735
respect (800)	0.720
cause (753)	0.704
beneficiary (715)	0.669
instrument (507)	0.650
time frame (45)	0.603
comparison (198)	0.600
frequency (78)	0.576
experiencer (259)	0.556
goal (282)	0.525
modality (17)	0.333
extent of space (67)	0.278
result (15)	0.222
condition (5)	0.000
property (6)	0.000

Automatic analysis: an example

– τὴν ἐνβολὴν καὶ σφόδρα ὑφ' ὑμῶν ἀμελουμένην ὄρω.

	FORM	LEMMA	POSTAG	HEAD	RELATION	SEMANTIC_ROLE
1	τὴν	ό	I-s---fa-	2	ATR	_
2	ἐνβολὴν	ἐνβολή	n-s---fa-	8	OBJ	stimulus
3	καὶ	καί	b-----	7	AuxY	_
4	σφόδρα	σφόδρα	d-----	7	ADV	degree
5	ὑφ'	ὑπό	r-----	7	AuxP	_
6	ὑμῶν	σύ	p-p---g-	5	OBJ	agent
7	ἀμελουμένην	ἀμελέω	v-sppefa-	2	ATR	_
8	όρω	όράω	v1spia---	0	PRED	_
9	.	.	u-----	0	AuxK	_