

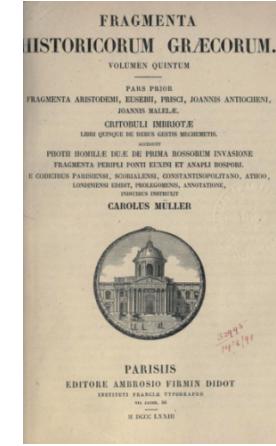
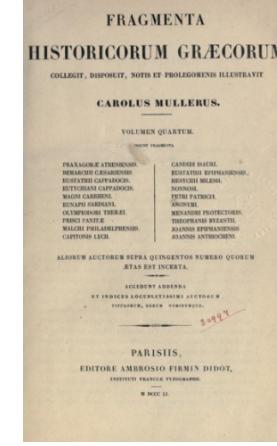
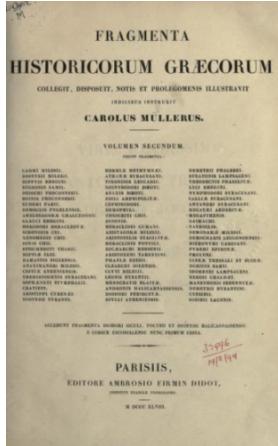
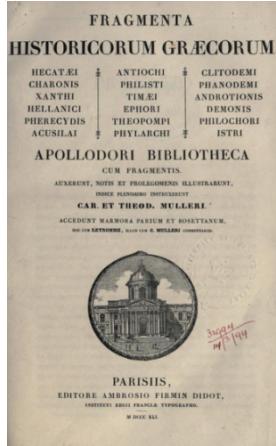


# The Digital Fragmenta Historicorum Graecorum

Tariq Yousef  
University of Leipzig

Sunoikisis DC 2016  
December 16, 2015  
Leipzig, Germany

# Digital Fragmenta Historicorum Graecorum (DFHG) Project



[www.dfhg-project.org](http://www.dfhg-project.org)

5 volumes - ca. 600 authors  
6th century BC - 7th century CE  
Greek-Latin alignment

# Dfhg-dev

Digital Fragmenta Historicorum Graecorum (DFHG): a machine-corrected version of the five volumes of Karl Müller's *Fragmenta Historicorum Graecorum*.

[View the Project on GitHub](#)  
OpenGreekAndLatin/dfhg-dev

Download  
**ZIP File**

Download  
**TAR Ball**

View On  
**GitHub**

This project is maintained by  
[OpenGreekAndLatin](#)

Hosted on GitHub Pages — Theme by [orderedlist](#)

# Welcome to dfhg-dev!

## About

This repository contains machine-corrected EpiDoc versions of the five volumes of Karl Müller's *Fragmenta Historicorum Graecorum* (FHG) as part of the [Digital Fragmenta Historicorum Graecorum \(DFG\)](#) project at the University of Leipzig. The files are encoded by the [Open Greek and Latin research team](#) (researchers and student assistants). While adequate, more can and should be done to improve the XML (the -dev(velopment) suffix indicates the repository is open to improvement).

If you want to help or contribute to the encoding process, please read the following documentation:

- [DFHG Guidelines \(1.0\)](#).
- [DFHG Authors](#).

For a full description of the DFHG project, please visit our [website](#).

## Contents

For a comprehensive view of the contents of the FHG, click [here](#). Authors encoded thus far are (authors are listed according to the order of Müller's FHG):

### Volume 1

- [Charon Lampsacenus](#)
- [Hellanicus Lesbius](#)
- [Acusilaus Argivus](#)
- [Antiochus Syracusanus](#)



GitHub, Inc. [US] [https://github.com/OpenGreekAndLatin/dfhg-dev/blob/master/Volume\\_1/Ister%20Cyrenaeus/dfhg1ister\\_grc.xml](https://github.com/OpenGreekAndLatin/dfhg-dev/blob/master/Volume_1/Ister%20Cyrenaeus/dfhg1ister_grc.xml)

```
132 <ab><title xml:lang="la">LIBER I.</title></ab>
133
134 <cit n="1" xml:lang="grc">
135   <bibl>Suidas v. Τιτανίδα γῆν</bibl>
136   <quote>... <persName role="author">Ιστρος</persName> ἐν <num value="1"
137     >α</num>
138     <title>Ἄττικῶν</title>, Τιτᾶνας βοῶν· ἔδ#976;οήθουν γὰρ τοῖς
139     ἀνθρώποις ἐπακούνοτες, ὡς <persName role="author"
140       >Νίκανδρος</persName> ἐν <num value="1">α</num>
141     <title>Αἰτωλικῶν</title>. </quote>
142     <note xml:lang="la">Idem legitur in <bibl>Photii Lexico</bibl>. </note>
143   </cit>
144
145   <cit n="2" xml:lang="grc">
146     <bibl>Apostolius XVIII , 77</bibl>
147     <quote>Τιτανίδα παροικεῖς, ἐπὶ τῶν φιλοθέων. Οἱ μὲν τὴν πᾶσαν γῆν, οἱ δὲ
148       τὴν Ἀττικὴν ἀπὸ Τιτηνίου, ἐνὸς τῶν Τιτάνων ἀρχαιοτέρου, οἰκήσαντος
149       περὶ Μαραθῶνα, δὲς μόνος οὐκ ἐστράτευσεν ἐπὶ τοὺς θεοὺς, ὡς <persName
150         role="author">Φιλόχορος</persName> ἐν <title>Τετραπόλει</title>
151         καὶ <persName role="author">Ιστρος</persName> ἐν <num type="ordinal"
152           value="1">πρώτῃ</num> τῶν <title>Ἄττικῶν</title>. </quote>
153     <note xml:lang="la">Vide <bibl>Pliavon. s. Τιταν. γῆν</bibl>, ubi eadem
154       leguntur, paucis omissis et suppresso Istri nomine. </note>
155   </cit>
156
157   <cit n="3" xml:lang="grc">
158     <bibl>Harpocrat.: Λαμπάς.</bibl>
159     <quote>... Τρεῖς ἄγουσιν Ἀθηναῖοι ἑορτὰς λαμπάδας (λαμπάδος <foreign
160       xml:lang="la">cod. Angl.</foreign>), Παναθηναίοις, καὶ
161       Ἑφαιστείοις, καὶ Προμηθείοις, ὡς <persName role="author"
162         >Πολέμων</persName> φησὶν ἐν τῷ περὶ τῶν ἐν τοῖς προπυλαιοῖς
163         πινάκων. <persName role="author">Ιστρος</persName> δ' ἐν <num
164           type="ordinal" value="1">πρώτῃ</num> τῶν <title type="alt"
165             >Ἀτθίδων</title>, εἰπὼν ὡς ἐν τῇ τῶν Απατουρίων ἑορτῇ Ἀθηναίων
166             οἱ καλλίστας στολὰς ἐνδεδυκότες, λαδ#976;όντες ἡμμένας λαμπάδας ἀπὸ
167             τῆς ἑστίας, ὑμοῦσι τὸν Ἑφαιστον θύοντες (<foreign xml:lang="la">sic
168               codd.</foreign>; <bibl>Valesius</bibl> θέοντες), ὑπόμνημα τοῦ
169             κατανοήσαντα τὴν χρείαν τοῦ πυρὸς διδάξαι τοὺς ἄλλους. </quote>
```

DFHG Authors - Progress

Last edit was made 2 days ago by Uta Koschmieder

File Edit View Insert Format Data Tools Add-ons Help

Arial 14 B I A DFHG pages Archive.org URL TLG names Fragmentary author canon. abbr. (cf. LSJ v.stoa.org/abbreviat XML encoding assigned

**FHG 1**

Praefatio  
De Vita et Scriptis Auctorum

Hecataei  
Charonis  
Xanthi  
Hellanic  
Pherecydis  
Acusilai

Apollodori Atheniensis Bibliothecae  
Fragmenta Bibliothecae

Antiochi  
Philisti

DFHG Authors

	A	B	C	D	E	F
1	DFHG Authors	DFHG pages	Archive.org URL	TLG names	Fragmentary author canon. abbr. (cf. LSJ v.stoa.org/abbreviat	XML encoding assigned
2						
3	<b>FHG 1</b>					
4	Praefatio	i-vii	<a href="https://archive.org/st">https://archive.org/st</a> -		-	MONICA BE
5	De Vita et Scriptis Auctorum	ix-xci	<a href="https://archive.org/st">https://archive.org/st</a> -		-	MONICA BE
6	Hecataei	1-31	<a href="https://archive.org/st">https://archive.org/st</a>	Hecataeus Milesius	Hecat	MONICA BE
7	Charonis	32-35	<a href="https://archive.org/st">https://archive.org/st</a>	Charon Lampsacenus	Charon	MONICA BE
8	Xanthi	36-44	<a href="https://archive.org/st">https://archive.org/st</a>	Xanthus Lydia	Xanth	MONICA BE
9	Hellanic	45-69	<a href="https://archive.org/st">https://archive.org/st</a>	Hellanicus of Lesbos	Hellanic	MONICA BE
10	Pherecydis	70-99	<a href="https://archive.org/st">https://archive.org/st</a>	Pherecydes Atheniensis	Pherecyd	MONICA BE
11	Acusilai	100-103	<a href="https://archive.org/st">https://archive.org/st</a>	Acusilaus Argivus	Acus	MONICA BE
12	Apollodori Atheniensis Bibliothecae	104-179	<a href="https://archive.org/st">https://archive.org/st</a>	Pseudo-Apollodorus (Scriptor Bibliothecae)	Ps-Apollod	AUTOMATIC
13	Fragmenta Bibliothecae	180	<a href="https://archive.org/st">https://archive.org/st</a>		FragBibl	AUTOMATIC
14	Antiochi	181-184	<a href="https://archive.org/st">https://archive.org/st</a>	Antiochus Syracusanus	AntiochHist	KEVIN STRAßBURG
15	Philisti	185-192	<a href="https://archive.org/st">https://archive.org/st</a>	Philistus Syracusanus	Philist	KEVIN STRAßBURG

DFHG Authors DFHG Editorial Team and Encoders

# Automatic XML Tagging



# Automatic XML Tagging

XML markup levels:

- **Macro Level Markup:** Chapters, sections, paragraph, ..etc
- **Micro Level Markup:** used to determine the data types of the elements (date, time, person name, ..etc) or the role of the elements (author, editor, work name, ..etc).

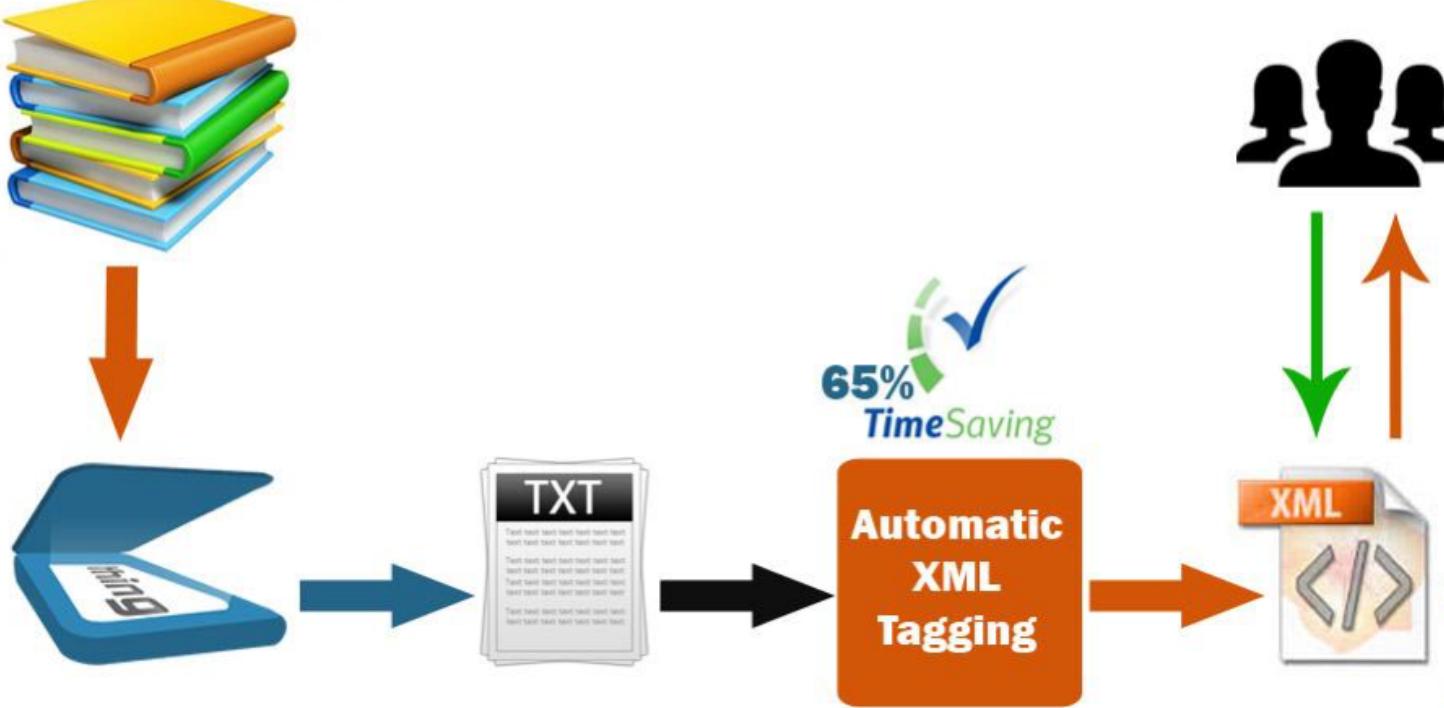
# DFHG Collection

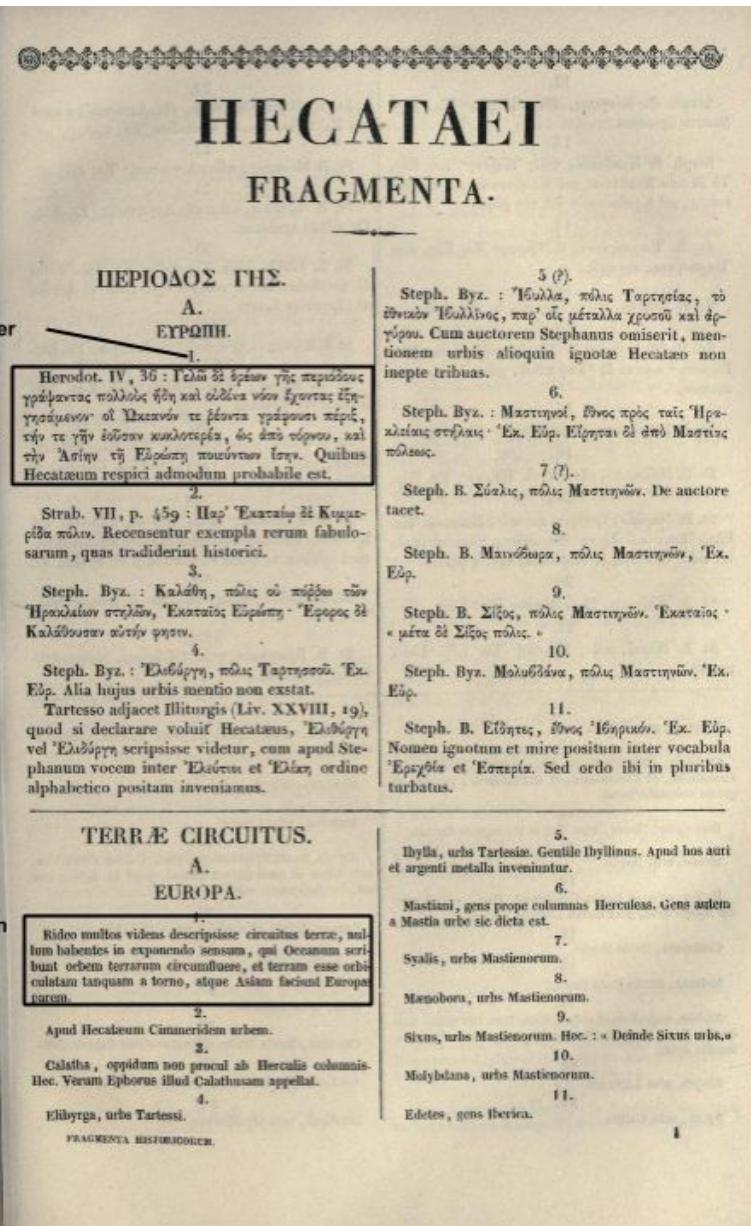
- 5 Volumes (available as image files)
- 600 Authors
- 3000 pages



- **Manual tagging is impractical in term of time and effort, especially when the data collection is huge.**

# Automatic XML Tagging





The plain texts were produced by an OCR.

### Greek Fragment (Original Image)

1.  
Herodot. IV, 36 : Γελῶ δὲ ὁρέων γῆς περιόδους γράψαντας πολλοὺς ἡδη καὶ οὐδένα νόον ἔχοντας ἐξηγησάμενον· οἱ Ωκεανόν τε ῥέοντα γράφουσι πέριξ, τὴν τε γῆν ἐσύσαν κυκλοτερέα, ὡς ἀπὸ τόρνου, καὶ τὴν Ἀσίην τῇ Εὐρώπῃ ποιεύντων ἴσην. Quibus Hecataeum respici admodum probabile est.

1.

<p>Herodot. IV, 36: Γελῶ δὲ ὁρέων γῆς περιόδους γράψαντας πολλοὺς ἡδη καὶ οὐδένα νόον ἔχοντας ἐξηγησάμενον οἱ Ωκεανόν τε ῥέοντα γράφουσι πέριξ, τὴν τε γῆν ἐσύσαν κυκλοτερέα, ὡς ἀπὸ τόρνου, καὶ τὴν Ασίην τῇ Εὐρώπῃ ποιεύντων ἴσην. Quibus Hecataeum respici admodum probabile est.</p>

```
<cit n="1" xml:lang="grc">
  <bibl>Herodot. IV, 36</bibl>
  <quote>Γελῶ δὲ ὁρέων γῆς περιόδους γράψαντας πολλοὺς ἡδη καὶ οὐδένα νόον ἔχοντας ἐξηγησάμενον· οἱ Ωκεανόν τε ῥέοντα γράφουσι πέριξ, τὴν τε γῆν ἐσύσαν κυκλοτερέα, ὡς ἀπὸ τόρνου, καὶ τὴν Ασίην τῇ Εὐρώπῃ ποιεύντων ἴσην.</quote>
  <note>Quibus Hecat&#230;um respici admodum probabile est. </note>
</cit>
```

XML Output

### Latin Fragment (Original Image)

1.  
Rideo multos videns descriptsse circuitus terrae, nullum habentes in exponendo sensum, qui Oceanum scribunt orbem terrarum circumfluere, et terram esse orbiculatam tanquam a torno, atque Asiam faciunt Europæ parem.

1.

<p>Rideo multos videns descriptsse circuitus terrae, nullum habentes in exponendo sensum, qui Oceanum scribunt orbem terrarum circumfluere, et terram esse orbiculatam tanquam a torno, atque Asiam faciunt Europæ parem.</p>

<p n="1">

Rideo multos videns descriptsse circuitus terrae, nullum habentes in exponendo sensum, qui Oceanum scribunt orbem terrarum circumfluere, et terram esse orbiculatam tanquam a torno, atque Asiam faciunt Europæ parem.

</p>

XML Output

# Fragment Boundary Identification

OCR output:

- <body> .... </body>
- <p>... </p>
- <pb n="287">

...  
<p>Ante Ogygen nihil apud Græcos narratur dignum quod  
memoretur, excepto Phoroneo, illius æquali, atque Inacho,  
Phoronei patre, qui primus Argolidis rex erat, ut narrat  
Acusilaus.  
14.  
<p>Acusilaus dicit Phoroneum fuisse primum hominum.  
... Hinc Plato in Timæo secutus Acusilaum, scribit:</p>  
</body>  
<pb n="101">

<body>

ex Codd. MSS. edita: Δευκαλίων, ἐφ' οὗ ὁ κατακλυσμὸς  
γέγονε, Προμηθέως μὲν ἦν νιὸς, μητρὸς δὲ ὡς  
πλεῖστοι λέγουσι Κλυμένης, ὡς δὲ Ήσιόδος Πρυνείης,  
ὡς δὲ Ἀκουσίλαος Ήσιόνης τῆς Νικεανοῦ καὶ τοῦ  
Προμηθέως.</p>

<p>Sturz. pro Ήσιόνης, quæ h. l. Oceanī filia dicitur,  
legendum putat Ασίας et pro corrupto  
illo Πρυνείης reponi vult Πανδώρας.</p>

<p>Schol. Pind. Ol. 9, 70: Κοινὰ τὰ περὶ Δευκαλίωνα  
καὶ Πύρφαν, καὶ δτὶ τοὺς λίθους κατόπιν βίπτοντες  
ἀνθρώπους ἐποίουν, μαρτυρεῖ Ἀκουσίλαος.</p>

8.

<p>Schol. Apollon. III, 1123: Ἀκουσίλαος καὶ Ήσιόδος  
ἐν ταῖς μεγάλαις Ήσιαῖς φασὶν ἔξ Ιοφώσσης τῆς  
Αἰήτου (Phrixum genuisse Argum).</p>

9.

<p>Schol. Apoll. Rh. IV, 1147: Περὶ δὲ τοῦ δέρους,  
δτὶ ἦν χρυσοῦν, οἱ πλεῖστοι ιστοροῦσιν. Ἀκουσίλαος δὲ  
ἐν τῷ περὶ γενεαλογιῶν πορφυρευθῆναι φησὶν ἀπὸ τῆς  
Θαλάσσης.</p>

10.

# Fragment Boundary Identification

2.

<p>Etymol. M.: Κοῖος· ὁ πατὴρ Λητοῦς· παρὰ τὸ κοεῖν, ὃ ἔστι νοεῖν καὶ συνιέναι, ὃ ἔστι συνετός. Κοῖον θ' Ὑπερίονά τ' Ἰαπετόν τε (ex Hesiodi Theog. 134). Οἱ Αἰολεῖς τῷ κ ἀντὶ τοῦ ν κέχρηνται. Οὗτοι δὲ Τιτᾶνες καὶ Τιτανίδες καλοῦνται, ὡς Ἀκουσίλαος.</p>

# Language Detection

2.

<p>Etymol. M.: Κοῖος· ὁ πατὴρ Λητοῦς· παρὰ τὸ κοεῖν, ὃ ἔστι νοεῖν καὶ συνιέναι, ὃ ἔστι συνετός. Κοῖον θ' Ὑπερίονά τ' Ἰαπετόν τε (ex Hesiodi Theog. 134). Οἱ Αἰολεῖς τῷ κ ἀντὶ τοῦ ν κέχρηνται. Οὔτοι δὲ Τιτᾶνες καὶ Τιτανίδες καλοῦνται, ὡς Ἀκουσίλαος.</p>

2.

<p>Cœus, pater Latonæ, sic nominatus a κοεῖν, quod idem est ac νοεῖν et συνιέναι (scire, intelligere), ita ut κοῖος sit quasi νοῖος prudens, intelligens. Nam Άeoles litera κ pro ν utuntur. Hesiodus: Cœum et Hyperiona et Iapetum. Hi vero Titanes et Titanides vocantur, ut ait Acusilaus.</p>

5 Latin words

36 Greek words

39 Latin words

6 Greek words

# Language Detection

2.

<p>Etymol. M.: Κοῖος· ὁ πατὴρ Λητοῦς· παρὰ τὸ κοεῖν, ὃ ἔστι νοεῖν καὶ συνιέναι, ὃ ἔστι συνετός. Κοῖον θ' Ὑπερίονά τ' Ἰαπετόν τε (ex Hesiodi Theog. 134). Οἱ Αἰολεῖς τῷ κ ἀντὶ τοῦ ν κέχρηνται. Οὗτοι δὲ Τιτᾶνες καὶ Τιτανίδες καλοῦνται, ὡς Ἀκουσίλαος.</p>

<cit n="2" xml:lang="grc">Etymol. M.: Κοῖος· ὁ πατὴρ Λητοῦς· παρὰ τὸ κοεῖν, ὃ ἔστι νοεῖν καὶ συνιέναι, ὃ ἔστι συνετός. Κοῖον θ' Ὑπερίονά τ' Ἰαπετόν τε (ex Hesiodi Theog. 134). Οἱ Αἰολεῖς τῷ κ ἀντὶ τοῦ ν κέχρηνται. Οὗτοι δὲ Τιτᾶνες καὶ Τιτανίδες καλοῦνται, ὡς Ἀκουσίλαος.</cit>

2.

<p>Cœus, pater Latonæ, sic nominatus a κοεῖν, quod idem est ac νοεῖν et συνιέναι (scire, intelligere), ita ut κοῖος sit quasi νοῖος prudens, intelligens. Nam Άοles litera κ pro ν utuntur. Hesiodus: Cœum et Hyperiona et Iapetum. Hi vero Titanes et Titanides vocantur, ut ait Acusilaus.</p>

<p n="2">Cœus, pater Latonæ, sic nominatus a κοεῖν, quod idem est ac νοεῖν et συνιέναι (scire, intelligere), ita ut κοῖος sit quasi νοῖος prudens, intelligens. Nam Άοles litera κ pro ν utuntur. Hesiodus: Cœum et Hyperiona et Iapetum. Hi vero Titanes et Titanides vocantur, ut ait Acusilaus.</p>

- We enclose every fragment in a <cit> element.
- We encode every fragment number by adding a @n attribute to the <cit> element.

# Fragment's Elements Identification

The Greek fragments usually contain the following elements:

- Fragment's number;
- Bibliographic reference to the source text that preserves the fragment;
- Greek text of the fragment (= quotation from the source text);
- Brief note in Lat

**OCR Output**

1.

<p>Herodot. IV, 36: Γελῶ δέ ὥρεον γῆς περιόδους  
γράψαντας πολλοὺς ἥρη καὶ σύδένα νόον ἔχοντας ἐξηγησάμενον  
οἱ Σκεανόν τε βέοντα γράφουσι πέρι,  
τὴν τα γῆν ἔσυσαν κυκλοτερέα, ὡς ἀπὸ τόρνου, καὶ  
τὴν Ασίην τῇ Εὐρώπῃ ποιεύντεν ἵστην. Quibus  
Hecataeum respici admodum probabile est.</p>

# Fragment's Elements Identification

- We enclose the bibliographic reference to the source text within a <bibl> element.
- The Greek text of the fragment is in <quote>.
- If the Greek contains a word or a sentence in Latin, we encode it with <foreign  
xml:lang="la">...</foreign>.
- Latin notes at the end of a fragment sit in <note> with an attribute @xml:lang="la".

# Fragment's Elements Identification

1.  
Herodot. IV, 36 : Γελῶ δὲ ὄρέων γῆς περιόδους  
γράψαντας πολλοὺς ἡδη καὶ οὐδένα νόον ἔχοντας ἐξη-  
γησάμενον· οἱ Ωκεανόν τε ῥέοντα γράφουσι πέριξ,  
τὴν τε γῆν ἐσῦσαν κυκλοτερέα, ὡς ἀπὸ τόρνου, καὶ  
τὴν Ἀσίην τῇ Εὐρώπῃ ποιεύντων ἵσην. Quibus  
Hecatæum respici admodum probabile est.

## OCR Output

1.  
<p>Herodot. IV, 36: Γελῶ δὲ ὄρέων γῆς περιόδους  
γράψαντας πολλοὺς ἡδη καὶ οὐδένα νόον ἔχοντας ἐξηγησάμενον·  
οἱ Ωκεανόν τε ῥέοντα γράφουσι πέριξ,  
τὴν τε γῆν ἐσῦσαν κυκλοτερέα, ὡς ἀπὸ τόρνου, καὶ  
τὴν Ἀσίην τῇ Εὐρώπῃ ποιεύντων ἵσην. Quibus  
Hecatæum respici admodum probabile est.</p>

<cit n="1" xml:lang="grc">  
<bibl>Herodot. IV, 36</bibl>  
<quote>Γελῶ δὲ ὄρέων γῆς περιόδους γράψαντας  
πολλοὺς ἡδη καὶ οὐδένα νόον ἔχοντας ἐξηγησάμενον· οἱ  
Ωκεανόν τε ῥέοντα γράφουσι πέριξ, τὴν τε γῆν ἐσῦσαν  
κυκλοτερέα, ὡς ἀπὸ τόρνου, καὶ τὴν Ἀσίην τῇ Εὐρώπῃ  
ποιεύντων ἵσην.</quote>  
<note>Quibus Hecat&#233;um respici admodum  
probabile est. </note>  
</cit>

## XML Output

[DEMO](#) [DEMO2](#)

# **The Greek - Latin Dynamic Lexicon**



# Goals

- Create a Greek-Latin lexicon using English as pivot language.
- Identify the translation relationships among the words in the parallel fragments of the DFHG bilingual corpus.
- Provide training data for a word alignment system.

Large digitized Ancient Greek-Latin lexica are publicly unavailable.

# Greek-Latin Dynamic Lexicon

**Data source(Perseus Digital Library):**

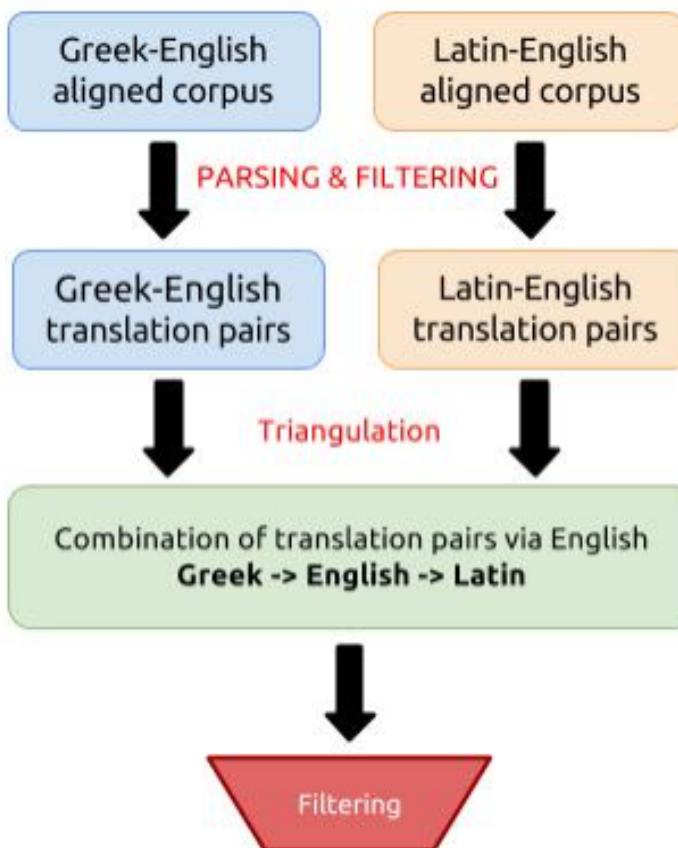
- Greek-English aligned texts (104 files)
- Latin-English aligned texts (59 files)



The parallel texts are aligned on a sentence level using **Moore's Bilingual Sentence Aligner**. Then the **Giza++** toolkit is used to align the sentence pairs at the level of individual words.

	Files	Sentences	Words	Distinct pairs
Ancient Greek	104	210 K	4,32 M	43 K
Latin	59	132 K	2,33 M	39 K

# The proposed method



# One-to-One Alignment

```
<sentence id="6">
  <wds lnum="L1">
    <w n="6-1" nrefs="6-3">illuc </w>
    <w n="6-2" nrefs="6-1">regredere</w>
    <w n="6-3" nrefs="6-4">ab</w>
    <w n="6-4" nrefs="6-6">ostio </w>
    <w n="6-5" nrefs="6-7">. </w>
  </wds>
  <wds lnum="L2">
    <w n="6-1" nrefs="6-2">get</w>
    <w n="6-2" nrefs="">away</w>
    <w n="6-3" nrefs="6-1">there </w>
    <w n="6-4" nrefs="6-3">from</w>
    <w n="6-5" nrefs="">the</w>
    <w n="6-6" nrefs="6-4">door</w>
    <w n="6-7" nrefs="6-5">!</w>
  </wds>
</sentence>
```

illuc regredere ab ostio .  
get away there from the door !

# One-to-many Alignment

```
<sentence id="8" >
  <wds lnum="L1">
    <w n="61-1" nrefs="8-1">καὶ</w>
    <w n="61-2" nrefs="8-2 8-3">πῶς</w>
    <w n="61-3" nrefs="8-5">λέγει</w>
    <w n="61-4" nrefs="8-6">;</w>
  </wds>
  <wds lnum="L2">
    <w n="61-1" nrefs="8-1">and</w>
    <w n="61-2" nrefs="8-2">what</w>
    <w n="61-3" nrefs="8-2">does</w>
    <w n="61-4" nrefs="">he</w>
    <w n="61-5" nrefs="8-3">say</w>
    <w n="61-6" nrefs="8-4">?</w>
  </wds>
</sentence>
```

Καὶ πῶς λέγει ;  
and what does he say ?

# Alignment errors

The Latin word (**et**) occurred 55034 times in the text and is aligned to:

Word	Frequency	Percentage
and	48185	88%
;	1422	2,5%
; and	1145	2%
,	818	1,5%
, and	545	1%
other translations	1919	4%

## Source of Errors:

The alignment is noisy and done automatically by the Giza++ software

# Alignment errors filtering

- Delete all entries that are punctuation or special characters on one side aligned to words on the other side.
- Delete translation pairs with probability less than 1%, because those pairs are noisy and incorrect.

Word	Frequency	Percentage
and	48185	88%
;	1422	2,5%
; and	1145	2%
,	919	1,5%
, and	545	1%
other translations	1919	4%

# Lemmatization

Latin and Greek are highly inflected languages, whereas English is considered a weakly inflected language.

eius cupio **filiam** virginem mihi desponderi.

i' d like his maiden **daughter** to be promised me in marriage.

piliae et **filiae** salutem.

love to pilia and your **daughter**.

mea haec erilis gestitavit **filia**.

this **daughter** of my mistress had them.

# Lemmatization

Latin and Greek are highly inflected languages, whereas English is considered a weakly inflected language.

eius cupio **filiam** virginem mihi desponderi.

i' d like his maiden **daughter** to be promised me in marriage.

piliae et **filiae** salutem.

love to pilia and your **daughter**.

mea haec erilis gestitavit **filia**.

this **daughter** of my mistress had them.

**Lemmatization** solves this problem

**filiam, filiae, filia** → **daughter**

**filia** → **daughter**

# Lemmatization

Latin and Greek are highly inflected languages, whereas English is considered a weakly inflected language.

καὶ ἐπειδὴ ἔτεκεν **uiόv**, ἔξαρνος ἦν μὴ εῖναι ἐξ αὐτοῦ τὸ παιδίον  
but when she gave birth to a **son**, (callias) denied that the child was his

**uiòv** δὲ Τεισίου τοῦ Ῥαμνουσίου  
the **son** of teisias of rhamnos

ούμὸς **uiόç**  
my **son**

**Lemmatization** solves this problem

**uiòv, uióç, uiόv** → **son**

**uiόç** → **son**

# MORPHEUS

Lemmatization of Latin and Greek words is done using **Morpheus**.

Morpheus is a morphological parsing and lemmatizing tool for Ancient Greek and Latin produced by the Perseus Project and available as a **web service**.

It takes as input a token (Greek or Latin) and returns the lemma from which this token derives, and the full morphological information (part of speech, case, number, gender, person, mood, etc.).

# Lemmatization

Lemmatization of English translations will produce better results.

E.g., translation candidates of the Greek word **λέγειν**

Translation	Freq	Precentage
say	551	36%
speak	492	32%
tell	149	9.7%
speaking	110	7%
said	89	6%
saying	54	3.5%
mention	45	2.9%
says	25	1.5%
spoke	19	1.2

Translation probabilities

Translation	Freq	Precentage
say	551	36%
speak	492	32%
tell	149	9.7%
<b> speak</b>	110	7%
<b> say</b>	89	6%
<b> say</b>	54	3.5%
mention	45	2.9%
<b> say</b>	25	1.5%
<b> speak</b>	19	1.2

Lemmatization of English words

Translation	Freq	Precentage
say	710	46.8%
speak	621	40.6%
tell	149	9.7%
mention	45	2.9%

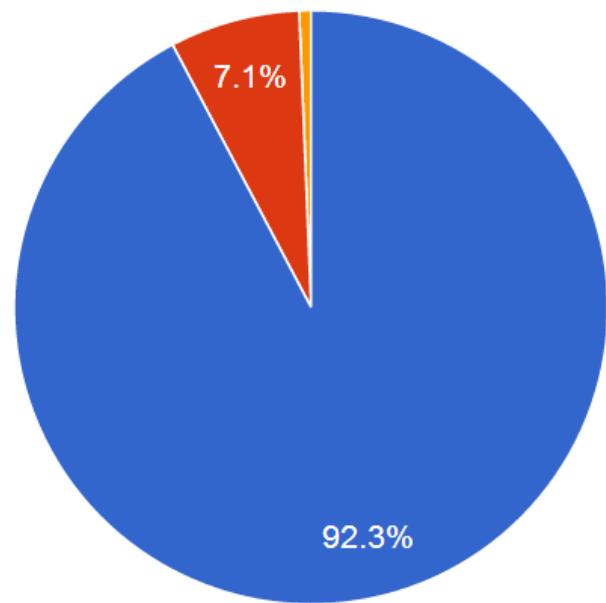
Group the results and recalculate  
the probabilities

# Triangulation

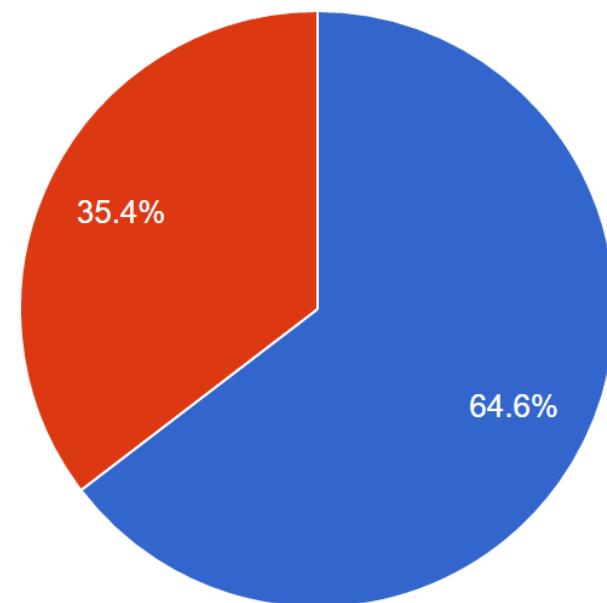
Triangulation is based on the assumption that two expressions are likely to be translations if they are translations of the same word in a third language.

## DEMO

Greek Translation of (city)



Latin Translation of (city)



- πόλις (πόλιν, πόλεως, πόλει, πόλιος)
- ἄστυ (ἄστυ, ἄστει, ἄστεως, ἄστεος)
- πολέω (πόλεις)

- civitas (civitate, civitatem, civitatis, civitas, civitati)
- urbs (urbem, urbe, urbis, urbs, urbi)

# Triangulation

- The English word (**city**) is translated to (92.3% **πόλις**), (7.1% **άστυ**), (0.6% **πολέω**)  
(**city**) is translated to (64.6% **civitas**), (35.4% **urbs**)
- The extracted pairs via triangulation:  
**(πόλις -civitas)**, **(πόλις -urbs)**, **(άστυ-civitas)**, **(άστυ-urbs)**, **(πολέω-civitas)**, **(πολέω-urbs)**.  
those pairs don't have the same level of relatedness,  
therefore we have to filter the results to keep only strong related pairs

# Translation-Pairs filtering

We used **Jaccard coefficient** as a similarity metric to measure the similarity or the relatedness between every Greek-Latin pairs.

$$J(A, B) = \frac{A \cap B}{A \cup B}$$

The relatedness between the Greek word (**πόλις**) and the Latin word (**civitas**) can be calculated as follows:

# Translation-Pairs filtering

πόλις		
city	7432	74 %
state	1911	19 %
town	366	3.5 %
athens	357	3.5 %

πολέω		
city	6673	64.3 %
of	2352	22.6 %
state	823	7.9 %
town	530	5.2 %

άστυ		
city	6747	91 %
citizen	399	4 %
town	236	3 %
others	96	1%

urbs		
city	44395	94.2 %
rome	2190	4.6 %

civitas		
city	8748	72.9 %
state	2340	19.5 %
citizenship	420	3.5 %
citizen	372	3.1 %

# Translation-Pairs filtering

<b>πόλις - civitas</b>	city, state	$(74 + 19 + 72 + 19.5)/200 = 92.25 \%$	92.25 %
<b>πόλις - urbs</b>	city, town	$(74 + 3.5 + 94.2 + 1.2)/200 = 86.45 \%$	86.45 %
<b>πολέω - civitas</b>	city, state	$(64.3 + 7.9 + 72 + 19.5)/200 = 81.85 \%$	81.85 %
<b>πολέω - urbs</b>	city, town	$(64.3 + 5.2 + 94.2 + 1.2)/200 = 82.45 \%$	82.45 %
<b>ἄστυ - civitas</b>	city, citizen	$(91 + 4 + 72 + 3.1)/200 = 85.05 \%$	85.05 %
<b>ἄστυ - urbs</b>	city, town	$(91 + 3 + 94.2 + 1.2)/200 = 94.7 \%$	94.7 %

# Translation-Pairs filtering

Ρώμη		
rome	8044	100 %

urbs		
city	44395	94.2 %
rome	2190	4,6 %
town	531	1.2 %

$$J(\text{ Ρώμη, urbs}) = (100+4.6)/200 = 51.3 \%$$



cor



PDF

## Translations

### Ancient Greek Translations

⌚ καρδία: κραδίη, καρδίᾳ, κραδίην

⌚ κῆρ: κῆρ

⌚ καρδιάω: καρδίαν, καρδία, καρδίας

⌚ ἡτορ: ἡτορ

⌚ θυμός: θυμός, θυμόν, θυμὸς, θυμῷ, θυμὸν

⌚ φρήν: φρένα, φρεσὶ, φρεσὶν, φρένας

⌚ ἡπαρ: ἡπαρ

↗ Morphological Analysis

🏷 Related Words

[cordi](#), [corde](#), [cordis](#), [corda](#), [mens](#)

## heart

- ✓ cor - κῆρ [κῆρ] (0,02)

ῆχθετο γὰρ κῆρ .

on this he groaned aloud ,and called upon his friend by name .

- ✓ cor - καρδίαν [καρδιάω] (0,04)

κάμε καρδίαν ἀμύσσει φροντίς .

my heart ,too ,is racked with anxiety ,and to you ,my friends ,will i make a disclosure .

- ✓ cor - κραδίην [καρδία] (0,07)

νηπύτι' ὡς ἄνοον κραδίην ἔχες .

Fool ,how witless is the heart thou hast !

- ✓ cor - καρδίᾳ [καρδία] (0,07)

τί ὅτι ἔθου ἐν τῇ καρδίᾳ σου τὸ πρᾶγμα τοῦτο ;

how is it that you have conceived this thing in your heart ?

- ✓ cor - κραδίῃ [καρδία] (0,07)

" τέτλαθι δή ,κραδίῃ .

" endure ,my heart ;

# Evaluation

Mean reciprocal rank (MRR)

$$MRR = \frac{1}{|Q|} \sum_{i=1}^{|Q|} RR_i$$

We selected randomly 200 translation pairs obtained via the proposed method with different frequencies (high and low) and different JACCARD Co values.

Each pair should be assigned into four categories: **Correct, small difference, big difference and incorrect.**

# Evaluation

Incorrect ✗	Big difference ⓘ	Small difference ⓘ	Correct ✈	πόλει - civitatis (city) [0.910379]
Incorrect ✗	Big difference ⓘ	Small difference ⓘ	Correct ✈	uiέ - gnato (son) [0.947044]
Incorrect ✗	Big difference ⓘ	Small difference ⓘ	Correct ✈	λέγοντες - diceret (say) [0.906456]
Incorrect ✗	Big difference ⓘ	Small difference ⓘ	Correct ✈	θεῶν - di (god) [0.935497]

Category	Reciprocal Rank
Correct	1
Small difference	0.75
Big difference	0.25
Uncorrect	0

Jaccard co	60% <	70% <	80% <	90% <
MMR	61.25 %	74 %	87.5 %	94.5 %

# Conclusion

The quality of the method depends on two factors:

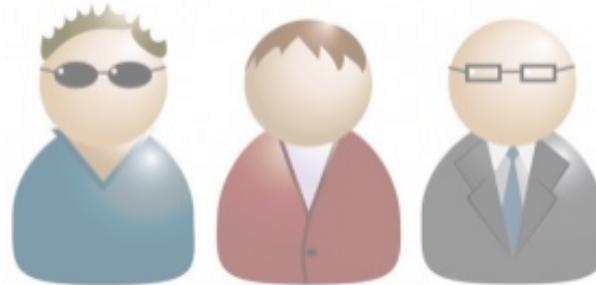
- **The size of aligned-parallel corpora**

Bigger corpora produce better translation probability distribution and more translation candidates.

- **The quality of the aligner**

Manually aligned corpora yield more accurate results.

# Named-Entities Recognition



# Named-Entities Recognition

**Named-Entity Recognition (NER)** is the task of **locating** and **classifying** single words and phrases in the text into predefined categories such as **person names**, **locations**, **organizations** as well as quantities, percentage, time expressions and monetary values.

Jim bought 300 shares of Acme Corp. in 2006.

[Jim]<sub>Person</sub> bought 300 shares of [Acme Corp.]<sub>Organization</sub> in [2006]<sub>Time</sub>.

# Entities Identification

Τὸ Σικελικὸν γένος ἔξελιπεν Ἰταλίαν· ως μὲν Ἐλλάνικος ὁ Λέσβιός φησι, τρίτη γενεᾶ πρότερον τῶν Τρωϊκῶν, Ἀλκυόνης ιεραμένης ἐν Ἀργει κατὰ τὸ ἔκτον καὶ εἰκοστὸν ἔτος. Δύο γὰρ ποιεῖ στόλους Ἰταλικοὺς διαβάντας εἰς Σικελίαν· τὸν μὲν πρότερον, Ἐλύμων, οὓς φησιν ὑπὸ Οἰνώτρων ἔχαναστηναι· τὸν δὲ μετὰ τοῦτον ἔτει πέμπτῳ γενόμενον, Αὔσονίων Ἰάπυγας φυγόντων· βασιλέα δὲ τούτων ἀπαφαίνει Σικελὸν, ἀφ' οὗ τοῦνομα τοῖς τε ἀνθρώποις καὶ τῇ νήσῳ τεθῆναι.

Genus **Siculum Italiam** reliquit, ut refert **Hellenicus Lesbius**, tribus ætatibus ante bellum Trojanum, anno vigesimo sexto sacerdotii quo **Alcyone Argis** exercebat. Duas enim classes **Italicas** in **Siciliam** trajecisse tradit: priorem **Elymorum**, quos ab Εnotris pulsos dicit: posteriorem vero **Ausonum Iapygas** fugientium, quæ quinto post anno eo trajecit. Horum autem regem **Siculum** fuisse dicit, a quo tam ipsi genti quam insulæ nomen inditum.

# Entities Identification

Τὸ Σικελικὸν γένος ἔξελιπεν Ἰταλίαν· ως μὲν Ἐλλάνικος ὁ Λέσβιός φησι, τρίτη γενεᾶ πρότερον τῶν Τρωϊκῶν, Ἀλκυόνης ἱεραμένης ἐν Ἀργει κατὰ τὸ ἔκτον καὶ εἰκοστὸν ἔτος. Δύο γὰρ ποιεῖ στόλους Ἰταλικοὺς διαβάντας εἰς Σικελίαν· τὸν μὲν πρότερον, Ἐλύμων, οὓς φησιν ὑπὸ Οἰνώτρων ἔχαναστηναι· τὸν δὲ μετὰ τοῦτον ἔτει πέμπτῳ γενόμενον, Αὔσονίων Ἰάπυγας φυγόντων· βασιλέα δὲ τούτων ἀπαφαίνει Σικελὸν, ἀφ' οὗ τοῦνομα τοῖς τε ἀνθρώποις καὶ τῇ νήσῳ τεθῆναι.

Genus **Siculum Italiam** reliquit, ut refert **Hellenicus Lesbius**, tribus ætatibus ante bellum Trojanum, anno vigesimo sexto sacerdotii quo **Alcyone Argis** exercebat. Duas enim classes **Italicas** in **Siciliam** trajecisse tradit: priorem **Elymorum**, quos ab Εnotris pulsos dicit: posteriorem vero **Ausonum Iapygas** fugientium, quæ quinto post anno eo trajecit. Horum autem regem **Siculum** fuisse dicit, a quo tam ipsi genti quam insulæ nomen inditum.

# Entities Identification

- All proper names are capitalized, but not all capitalized words are proper names
- The POS of a NE must be **Noun** or **UNKNOWN**
- Since the Latin fragments are translations of the Ancient Greek fragments, each proper name in the Ancient Greek fragment has a correspondence in the Latin parallel fragment, and they are similar from a syntactic or phonetic point of view.

# Entities Identification

In order to calculate the similarity between the translation pairs, we need a metric to measure the similarity. **Levenshtein Edit Distance** is the most popular metric

F1: Token begins with a capital letter

F2: POS tag

F3: Levenshtein distance

Words pair	F1	F2	F3
Αθήνας - Athenas	1	Noun - Noun	0
Άργος - Argos	1	Noun - Noun	0
Έλλανικος - Hellanico	1	Noun / Unknown	1

# Transliteration

converting the text from one writing system to another.

GRC	EN	GRC	EN	GRC	EN
A, α	A, a	N, ν	N, n	I, ι	I, i
B, β	B, b	Ξ, ξ	X, x	K, κ	K, k
Γ, γ	G, g	Ο, ο	O, o	Λ, λ	L, l
Δ, δ	D, d	Π, π	P, p	M, μ	M, m
Ε, ε	E, e	Ρ, ρ	R, r	Φ, φ	Ph, ph
Ζ, ζ	Z, z	Σ, σ	S, s	X, x	Ch, ch
Η, η	Ē, ē	Τ, τ	T, t	Ψ, ψ	Ps, ps
Θ, θ	Th, th	Υ, υ	Y, y	Ω, ω	Ō, ō

Source Text (Ancient Greek)

Τύραητα, πόλις Σαυνιτῶν. Φίλιστος ἐνδεκάτῳ.

Latin Transliteration

Týraēta, pólis Saunitōn. Phílistos hendekátōi.

Transliteration (without accent) **Tyraeta**, polis **Sauniton**. **Philistos** hendekatoi.

Target Text (Latin)

**Tyrseta**, urbs **Samnitum**. **Philistus** undecimo.

# Levenshtein Distance

Levenshtein distance is the most common metric for measuring the distance between two words or strings in terms of minimum number of insertion, deletion and substitution required to change one word or sequence into the other

$$\text{lev}_{a,b}(i,j) = \begin{cases} \max(i,j) \\ \min \begin{cases} \text{lev}_{a,b}(i-1, j) + 1 \\ \text{lev}_{a,b}(i, j-1) + 1 \\ \text{lev}_{a,b}(i-1, j-1) + 1_{(a_i \neq b_j)} \end{cases} \end{cases}$$

Ancient Greek	Latinized Greek	Latin	Levenshtein Distance
Φίλιστος	Philistos	Philistus	1
Έλλανικος	Hellanikos	Hellanicus	2
Δαμασκὸν	Damaskos	Damascus	2
Κύπρος	kypros	Cyprus	2

# Modified Levenshtein Distance

Levenshtein distance is not very helpful in our case, because it is binary and there is no tolerance with transliteration errors.

the distance between letters is not binary, but it is on scale. The cost of insertion or deletion depends on:

- Letter position (part of the prefix or suffix)
- Letter type (vowel or consonant)

Group elements	Distance within group elements
E, e Ē, ē I, i Y, y	0.15
Ō, ū O, o U, u Y, y	0.15
C, c K, k Q, q	0.15
S, s C, c	0.25
M, m N, n	0.35

# Modified Levenshtein Distance

Ancient Greek	Latinized Greek	Latin	Levenshtein Distance	Modified Levenshtein	<i>Dist/Length</i>
Φίλιστος	Philistos	Philistus	1	0.15	0,0167
Έλλανικος	Hellanikos	Hellenicus	2	0.30	0.03
Δαμασκόν	Damaskos	Damascus	2	0.30	0.0375
Κύπρος	kypros	Cyprus	2	0.30	0.05

# Relations between Entities

Τὸ Σικελικὸν γένος ἔξελιπεν Ἰταλίαν· ως μὲν Ἐλλάνικος ὁ Λέσβιός φησι, τρίτη γενεᾶ πρότερον τῶν Τρωϊκῶν, Ἀλκυόνης ιεραμένης ἐν Ἀργει κατὰ τὸ ἔκτον καὶ εἰκοστὸν ἔτος. Δύο γὰρ ποιεῖ στόλους Ἰταλικοὺς διαβάντας εἰς Σικελίαν· τὸν μὲν πρότερον, Ἐλύμων, οὓς φησιν ὑπὸ Οἰνώτρων ἔχαναστηναι· τὸν δὲ μετὰ τοῦτον ἔτει πέμπτῳ γενόμενον, Αὔσονίων Ἰάπυγας φυγόντων· βασιλέα δὲ τούτων ἀπαφαίνει Σικελὸν, ἀφ' οὗ τοῦνομα τοῖς τε ἀνθρώποις καὶ τῇ νήσῳ τεθῆναι.

Genus **Siculum Italiam** reliquit, ut refert **Hellenicus Lesbius**, tribus ætatibus ante bellum Trojanum, anno vigesimo sexto sacerdotii quo **Alcyone Argis** exercebat. Duas enim classes **Italicas** in **Siciliam** trajecisse tradit: priorem **Elymorum**, quos ab Εnotris pulsos dicit: posteriorem vero **Ausonum Iapygas** fugientium, quæ quinto post anno eo trajecit. Horum autem regem **Siculum** fuisse dicit, a quo tam ipsi genti quam insulæ nomen inditum.

# Relations between Entities

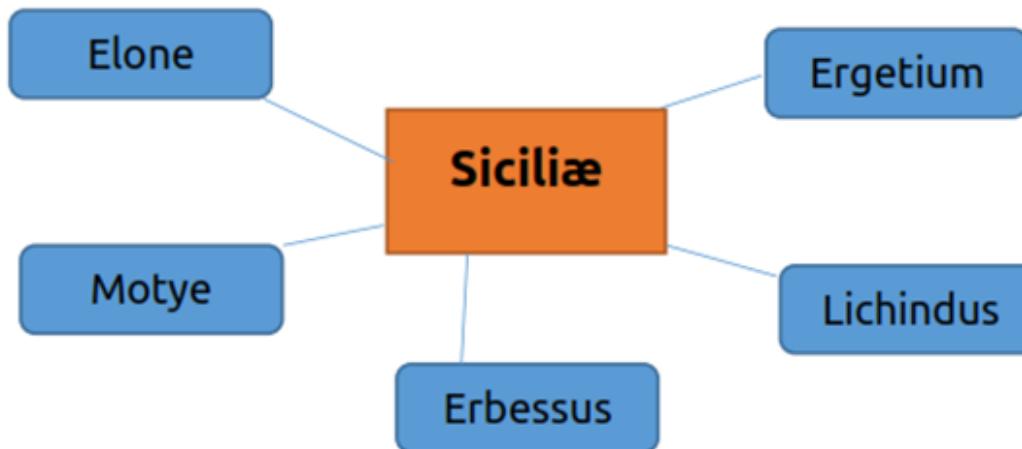
Phrase	Frequency in the text	Category
Hellenicus Lesbius	6	Author Name
Charon Lampsacenus	5	Author Name
Antiochus Syracusanus	4	Author Name
Acusilaus Argivus	1	Author Name

# Relations between Entities

Phrase	First NE	Relation	Second NE
Apollinis <b>ex</b> Callisto	Apollinis	<b>ex</b>	Callisto
Phœnissas <b>ex</b> Cypro	Phœnissas	<b>ex</b>	Cypro
Ergetium, <b>urbs</b> Siciliæ	Ergetium	, <b>urbs</b>	Siciliæ
Lichindus, <b>urbs</b> Siciliæ	Lichindus	, <b>urbs</b>	Siciliæ
Erbessus, <b>urbs</b> Siciliæ	Erbessus	, <b>urbs</b>	Siciliæ
Motye, <b>urbs</b> Siciliæ	Motye	, <b>urbs</b>	Siciliæ
Elone, <b>urbs</b> Siciliæ	Elone	, <b>urbs</b>	Siciliæ
Hellenicus <b>libro primo</b> Deucalioniae	Hellenicus	<b>libro primo</b>	Deucalioniae
Hellenico <b>libro primo</b> Lesbiorum	Hellenico	<b>libro primo</b>	Lesbiorum
Hellenico <b>libro primo</b> Troicorum	Hellenico	<b>libro primo</b>	Troicorum
Marsyas <b>libro primo</b> Macedonicorum	Marsyas	<b>libro primo</b>	Macedonicorum
Clidemus <b>libro primo</b> Protagoniæ	Clidemus	<b>libro primo</b>	Protagoniæ
Clidemus <b>libro primo</b> Atthidis	Clidemus	<b>libro primo</b>	Atthidis

# Relations between Entities

Ergetium, <b>urbs</b> Siciliæ	Ergetium	, urbs	Siciliæ
Lichindus, <b>urbs</b> Siciliæ	Lichindus	, urbs	Siciliæ
Erbessus, <b>urbs</b> Siciliæ	Erbessus	, urbs	Siciliæ
Motye, <b>urbs</b> Siciliæ	Motye	, urbs	Siciliæ
Elone, <b>urbs</b> Siciliæ	Elone	, urbs	Siciliæ



# Named-Entities Classification

**PLEIADES** is a gazetteer containing information about ancient world places and their locations and names.

<http://pleiades.stoa.org/>

# Named-Entities Classification

*A Dictionary of Greek and Roman Antiquities (1890)*

William Smith, LLD, William Wayte, G. E. Marindin, Ed

<http://www.perseus.tufts.edu/hopper/text?doc=Perseus%3atext%3a1999.04.0063>

## GeoNames

The GeoNames geographical database covers all countries and contains over eight million placenames that are available for download free of charge.

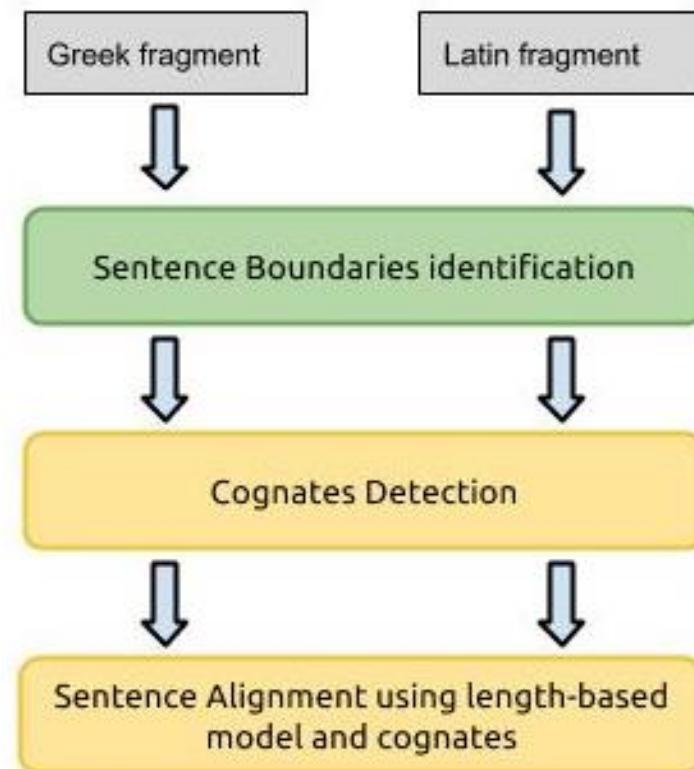
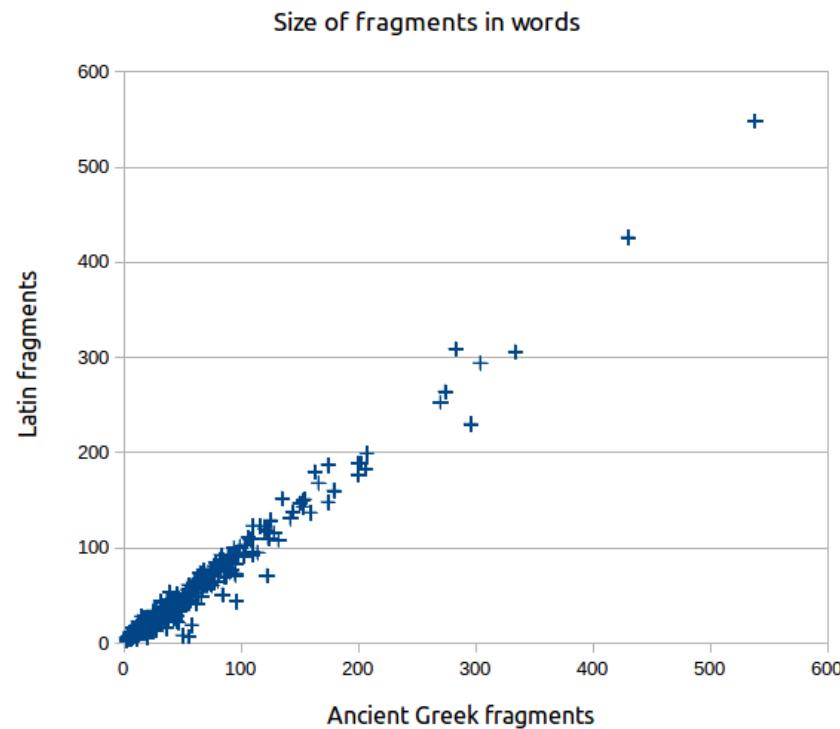
<http://www.geonames.org/>

# Word Alignment



# Word Alignment

1. Sentence Boundary Identification
2. Sentence Alignment



# Word Alignment model

**Μακτώριον, πόλις Σικελίας· Φίλιστος πρώτω, ἥν ἔκτισε μόνην. Τὸ ἔθνικὸν, Μακτωρῖνος.**

**Mactorium, urbs Siciliæ, de qua Philistus libro primo; quam condidit solam. Gentile, Mactorinus.**

## 1. Entities Alignment (*Transliteration & Modified Levenshtein distance*)

## 2. Dynamic Lexicon pairs (*Jaccard coefficient*)

## 2. Dynamic Lexicon pairs (*Jaccard coefficient*)

### 3. Co-occurrence Significance

Biemann and Quasthoff, 2005

Given two words A (Greek), B (Latin), each occurring a, b times in sentences (respectively), and k times together. We calculate the significance  $\text{sig}(A, B)$  of their occurrence (trans-co-occurrence) as follows:

For instance, (**πόλις**) occurs 55 times in the Greek fragments, and (**urbs**) occurs 57 times in the Latin fragments, and they occur 45 times

$$\text{sig}(A, B) = \frac{x - k \log x + \log k!}{\log n}$$

with       $x = \frac{ab}{n}$       , n: number of sentences

19,36

Sig (**πόλις, urbs**) =

	libro	primo	quam	condidit	solam	Gentile	Mactorinus
<b>ἐθνικὸν</b>	1.58	0.15	0.14	0.38	0.63	<b>6.39</b>	0.61

## 3.

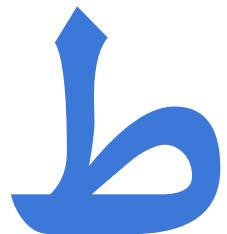
**Co-occurrence Significance**

Biemann and Quasthoff, 2005

	Μακτώριον	πόλις	Σικελίας	Φίλιστος	πρώτω	ἢν	ἔκτισε	μόνην	Τὸ	έθνικὸν	Μακτωρῖνος
Mactorium											
urbs											
Siciliæ											
de											
qua											
Philistus											
libro							0.15			1.58	
primo						0.18				0.15	
quam						1.04				0.14	
condidit						2.85				0.38	
solam						0.71				0.63	
Gentile						0.38			6.39		
Mactorinus						0.71				0.61	

# Morphological Disambiguator

[Demo](#)



THANKS!