

Corpus-driven Greek Language Learning

James K. Tauber

jktauber.com

[@jttauber](https://twitter.com/jttauber)

C. S. Lewis 1/2

“The great gain was that I very soon became able to understand a great deal without (even mentally) translating it; I was beginning to think in Greek.

That is the great Rubicon to cross in learning any language.

Those in whom the Greek word lives only while they are hunting for it in the lexicon, and who then substitute the English word for it, are not reading the Greek at all; they are only solving a puzzle.

C. S. Lewis 2/2

“The very formula, ‘Naus means a ship’, is wrong. Naus and ship both mean a thing, they do not mean one another.

Behind Naus...we want to have a picture of a dark, slender mass with sail or oars, climbing the ridges, with no officious English word intruding.”

Surprised by Joy: The Shape of My Early Life, pp. 140–141.

The Solution

“... comprehensible input is the crucial and necessary ingredient for the acquisition of language.”

Stephen Krashen

The Traditional Approach

vocabulary driven by paradigm being learnt

vocabulary not shown in context

hard to show much real text early on

The Myth of Vocabulary Coverage for the Greek New Testament

The **10** most common words account
for **37%** of the text

The **100** most common words account
for **66%** of the text

The Myth of Vocabulary Coverage for the Greek New Testament

If you learn the **100 most common words**,
you'll...

- know at least one word in **99.9%** of verses
- know at least **50%** of words in **91.3%** of verses
- know at least **75%** of words in **24.4%** of verses
- know at least **90%** of words in **2.1%** of verses
- know at least **95%** of words in **0.6%** of verses
- know **all** words in **0.4%** of verses

Levels of Understanding



too easy for the student

the right level for **speed reading**

the right level for **extensive reading**

the right level for **intensive reading**

the right level for **intensive reading**
with suitable scaffolding

too hard or requiring too much scaffolding

Use Cases

warn teacher / content author where
passage might need adaptation

search for passages appropriate for a
particular student

sort passages in a graded sequence

Example at 98%

You live and work in Tokyo. Tokyo is a big city. More than 13 million people live around you. Of course you are never, but you are always lonely.

Every morning, you get up and take the train to work. Every night, you take the train again to go home. The train is always crowded.

When people ask about your work, you tell them, “I move papers around.” It’s a joke, but it’s also true. You don’t like your work.

Tonight you are returning home. It’s late at night. Is sleeping no one. Sometimes you see don’t to sleep all day.

You are tired. You are so tired...

<https://magisterp.com/2016/08/21/how-comprehensible-must-reading-be/>

Example at 95%

In the morning, you start again. You shower, get dressed, and for pocklent walk. You move slowly, half-awake. Then, suddenly, you stop. Something is different. Are streets of fossit. Really of fossit.

There are no people. No cars. Nothing. “A where dowargle?” you ask yourself.

Suddenly, there is by quapen loud—a police car. It speeds by and almost hits you. It crashes into a store across the street!

Then, another for farfoofles police car. The police officer sees you. “Off the street!” he shouts. “Go home, lock your door!”

“What? Why?” you shout back. But it’s too late. He is gone.

<https://magisterp.com/2016/08/21/how-comprehensible-must-reading-be/>

Example at 90%

“Of prippy fy what’s?” you ask yourself. Suddenly, a man runs by. Viggling toward he is crawn for kofoon. There is blood all over his shirt.

“In order to baboot!” you shout, but he doesn’t stop. You follow him. For kofoon outside, you stop. Is lying loopity of on the ground. She is not moving.

“Hey!” you shout. “Are you OK?” She doesn’t answer. Are closed her nawiesly. Her fingers chay are moving. Open, close; open, close.

“She’s alive!” you say to yourself. “No! Her don’t of gleep!” someone the frickles. You look up. Three people are waving at you from across the street.

<https://magisterp.com/2016/08/21/how-comprehensible-must-reading-be/>

Example at 80%

“Of bingle for help!” you shout. “This is dying loopitish!” You put your fingers on her neck. Nothing. Flid her not a weafling is.

You take out by joople your and 119 the of bingle, the emergency number in Japan. There’s no answer!

Then muchy you that you have befourn assengle a new. It’s with gutring your, Evie. Hunwres she at Tokyo University. You assengle play.

“...if you get this...” Evie says. “... now vickarn I can’t... of passit important is...” Suddenly, to dingle, around looks she. “Oh no, they’re here! Cripettly... of frib!

Them wasple OF FRIB ON!...” BEEP! of the assengle to parantles. Then gratoonly something behind you...

<https://magisterp.com/2016/08/21/how-comprehensible-must-reading-be/>

The Myth of Vocabulary Coverage for the Greek New Testament

If you learn the **100 most common words**,
you'll...

- know at least one word in **99.9%** of verses
- know at least **50%** of words in **91.3%** of verses
- know at least **75%** of words in **24.4%** of verses
- know at least **90%** of words in **2.1%** of verses
- know at least **95%** of words in **0.6%** of verses
- know **all** words in **0.4%** of verses

Frequency

Frequency-based vocabulary lists are important post-beginner.

For beginners, passages either need to be constructed, heavily adapted, or we must not rely on corpus frequency as the primary ordering priority.

Pericope Experiment

1. 1 John 4:1–4:21
2. 1 John 5:1–5:12 (80.31%; 25 additional words)
3. 2 John 1:1–1:3 (84.75%; 8 additional words)
- ...
11. John 12:42–12:50 (92.73%; 12 additional words)
12. John 1:1–1:9 (95.33%; 5 missing words)
13. John 5:31–5:47 (93.49%)

Caveats and Limitations

Didn't consider:

- necessary spaced repetition
- inflectional morphology
- syntax
- idioms
- derivational morphology / cognates

Why weren't they considered?

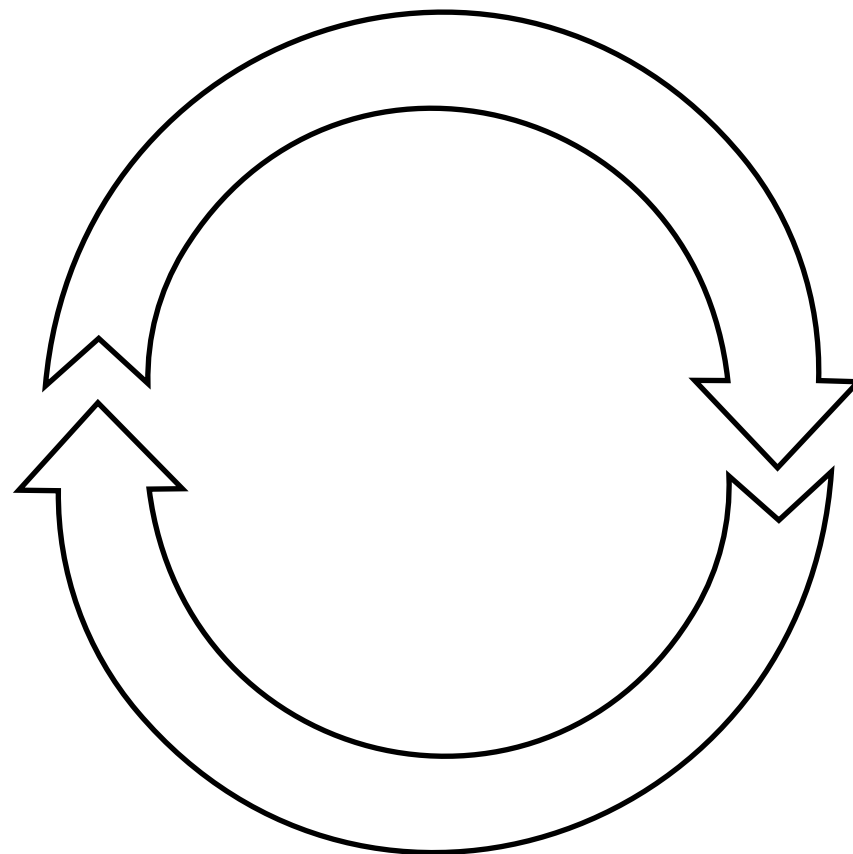
Because they weren't in an open,
machine-actionable format!

Linguistic Databases

Motivate

Enable

Learning Tools



Research Program

*linguistics,
descriptive grammar*

*digital philology,
corpus linguistics*

*learning science,
applied linguistics*

how to model
language

how to model
text

how to model
knowledge

*educational statistics,
learning analytics,
data science*

MorphGNT

041801	RD	----	APN-	Ταῦτα Ταῦτα ταῦτα οὗτος
041801	V-	-	AAPNSM-	ἔειπὼν εἰπὼν εἰπὼν λέγω
041801	N-	----	NSM-	Ἰησοῦς Ἰησοῦς Ἰησοῦς Ἰησοῦς
041801	V-	3AAI-	S--	ἐξῆλθεν ἐξῆλθεν ἐξῆλθε(ν) ἐξέρχομαι
041801	P-	-----		σὺν σὺν σύν σύν
041801	RA	----	DPM-	τοῖς τοῖς τοῖς ὁ
041801	N-	----	DPM-	μαθηταῖς μαθηταῖς μαθηταῖς μαθητής
041801	RP	----	GSM-	αὐτοῦ αὐτοῦ αὐτοῦ αὐτός
041801	P-	-----		πέραν πέραν πέραν πέραν

Word Lists by Frequency

```
awk '{print $7}' | sort |  
    uniq -c | sort -r
```

Analytical Lexicon

```
awk '{print $6,$2,$3,$7}' |  
    sort | uniq
```



James Tauber

@jtauber



Those of you who know me well understand why every new Greek Reader published makes me simultaneously jump for joy and cry a little.

12:05 AM - 13 Oct 2018

1 Retweet 13 Likes



2 1 13



James Tauber @jtauber · 13 Oct 2018



jump for joy because I'm a huge believer in comprehensible input with appropriate scaffolding and respect tremendously the effort that goes into producing these readers

1 2



James Tauber @jtauber · 13 Oct 2018



but cry a little because, as I've ranted for years, it's such a shame this data is locked in printed books rather than published in an open, machine actionable format with the printed version just generated from that data

1 1 11

Readers

4.2.1 Καὶ τῇ ἡμέρᾳ τῇ τρίτῃ γάμος¹ ἐγένετο ἐν Κανὰ² τῆς Γαλιλαίας, καὶ ἦν ἡ μήτηρ τοῦ Ἰησοῦ ἐκεῖ. **2** ἐκλήθη δὲ καὶ ὁ Ἰησοῦς καὶ οἱ μαθηταὶ αὐτοῦ εἰς τὸν γάμον.³ **3** καὶ ὑστερήσαντος⁴ οἴνου λέγει ἡ μήτηρ τοῦ Ἰησοῦ πρὸς αὐτόν· Οἶνον οὐκ ἔχουσιν. **4** καὶ λέγει αὐτῇ ὁ Ἰησοῦς· Τί ἐμοὶ καὶ σοί, γύναι; οὕτω⁵ ἦκει⁶ ἡ ὥρα μου. **5** λέγει ἡ μήτηρ αὐτοῦ τοῖς διακόνοις· Ὅτι ἂν λέγῃ ὑμῖν ποιήσατε. **6** ἦσαν δὲ ἐκεῖ λίθιναι⁸ ὑδρίαι⁹ ἑξ¹⁰ κατὰ τὸν καθαρισμὸν¹¹ τῶν Ἰουδαίων κείμεναι,¹² χωροῦσαι¹³ ἀνὰ¹⁴ μετρητὰς¹⁵ δύο ἢ τρεῖς. **7** λέγει αὐτοῖς ὁ Ἰησοῦς· Γεμίσατε¹⁶ τὰς ὑδρίας¹⁷ ὕδατος· καὶ ἐγέμισαν¹⁸ αὐτὰς ἕως ἄνω.¹⁹ **8** καὶ λέγει αὐτοῖς· Ἀντλήσατε²⁰ νῦν καὶ φέρετε τῷ ἀρχιτρικλίνῳ.²¹ οἱ δὲ ἤνεγκαν. **9** ὡς δὲ ἐγεύσατο²² ὁ ἀρχιτρίκλινος²³ τὸ ὕδωρ οἶνον γεγεννημένον, καὶ οὐκ ᾔδει πόθεν²⁴ ἐστίν, οἱ δὲ διάκονοι²⁵ ᾔδειςαν οἱ ἡντληκότες²⁶ τὸ ὕδωρ, φωνεῖ τὸν νυμφίον²⁷ ὁ ἀρχιτρίκλινος²⁸ **10** καὶ λέγει αὐτῷ· Πᾶς ἄνθρωπος πρῶτον τὸν καλὸν οἶνον τίθησιν, καὶ ὅταν μεθυσθῶσιν²⁹ τὸν ἐλάσσω.³⁰ σὺ τετήρηκας τὸν καλὸν οἶνον ἕως ἄρτι. **11** ταύτην ἐποίησεν ἀρχὴν τῶν σημείων ὁ Ἰησοῦς ἐν Κανὰ³¹ τῆς Γαλιλαίας καὶ ἐφάνέρωσεν τὴν δόξαν αὐτοῦ, καὶ ἐπίστευσαν εἰς αὐτὸν οἱ μαθηταὶ αὐτοῦ.

Readers

¹γάμος, ου, ό – a marriage, wedding, wedding-feast

²Κανά, ή – Cana

³γάμος, ου, ό – a marriage, wedding, wedding-feast

⁴ύστερέω – AAP GSM – I am lacking, fall short, suffer need

⁵οὔπω – not yet

⁶ἤκω – PAI 3S – I have come, am present

⁷διάκονος, ου, ό/ή – a waiter, servant, administrator

⁸λίθινος, η, ον – made of stone

⁹ύδρία, ας, ή – a water pot

¹⁰ἕξ – six

¹¹καθαρισμός, οὔ, ό – cleansing, purifying, purification

¹²κεῖμαι – PMP NPF – I lie, recline, am laid

¹³χωρέω – PAP NPF – I make room, go, receive

¹⁴ἀνά – and, apiece, by, each, every, in, through

¹⁵μετρητής, οὔ, ό – a measure, amphora

¹⁶γεμίζω – AAD 2P – I fill, load

¹⁷ύδρία, ας, ή – a water pot

¹⁸γεμίζω – AAI 3P – I fill, load

¹⁹ἄνω – up, above, things above, heaven

²⁰ἀντλέω – AAD 2P – I draw, draw out

²¹ἀρχιτρίκλινος, ου, ό – master of the feast

²²γεύομαι – AMI 3S – I taste, experience

²³ἀρχιτρίκλινος, ου, ό – master of the feast

²⁴πόθεν – whence

²⁵διάκονος, ου, ό/ή – a waiter, servant, administrator

²⁶ἀντλέω – XAP NPM – I draw, draw out

²⁷νυμφίος, ου, ό – a bridegroom

²⁸ἀρχιτρίκλινος, ου, ό – master of the feast

²⁹μεθύω – APS 3P – I am drunk

³⁰ἐλάσσων/ἐλασσον – less, smaller, inferior

³¹Κανά, ή – Cana

Generating a Greek Reader

```
./frequency_exclusion.py 31 > exclude.txt
```

```
# edit exclude to your liking
```

```
./make_glosses.py "John 4:1-11"  
  --exclude exclude.txt  
  > glosses.yaml
```

```
# edit glosses.yaml to your liking
```

```
./make_headwords.py "John 4:1-11"  
  --exclude exclude.txt  
  > headwords.yaml
```

Generating a Greek Reader

```
./reader.py "John 4:1-11"  
  --headwords headwords.yaml  
  --glosses glosses.yaml  
  --exclude exclude.txt  
  --typeface "Skolar PE"  
> reader.tex
```

Generating a Greek Reader

`https://github.com/
jtauber/greek-reader`

Written in Python 3 and open
source under an MIT license

Generating a Greek Reader

- electronic text (*SBLGNT*)
- lemmatisation (*MorphGNT*)
- parse codes (*MorphGNT*)
- glosses (*Dodson*)

Books become “UI”
derived from databases

Morphology

inflected forms versus lemmas

lists versus rules

lexical redundancy

Cost of Learning a Form

λέγει > λέγεις

λέγει > εἶπεν

λέγει > ἀντιλέγει ?

λέγει > λόγος ?

Ordering Within Morphology

no need to wait to learn about
athematic verbs and **δίδωμι** to
learn **ἔδωκεν** or **δός**.

Ordering Within Morphology

ἔδωκεν

δοῦναι

ἐδόθη

δώσω

δώσει

Ten Most Common Verb Parses (out of 379)

aorist active 3rd singular

present active 3rd singular

aorist active 3rd plural

aorist active infinitive

present active participle nominative singular masculine

aorist active participle nominative singular masculine

imperfect active 3rd singular

present active 1st singular

present active infinitive

present active participle nominative plural masculine

Morphological Lexicon

principal parts

stem relationships

word formation

transparency of lexical relatedness

Lexical Relatedness

Ἰταλία:Ἰταλικός

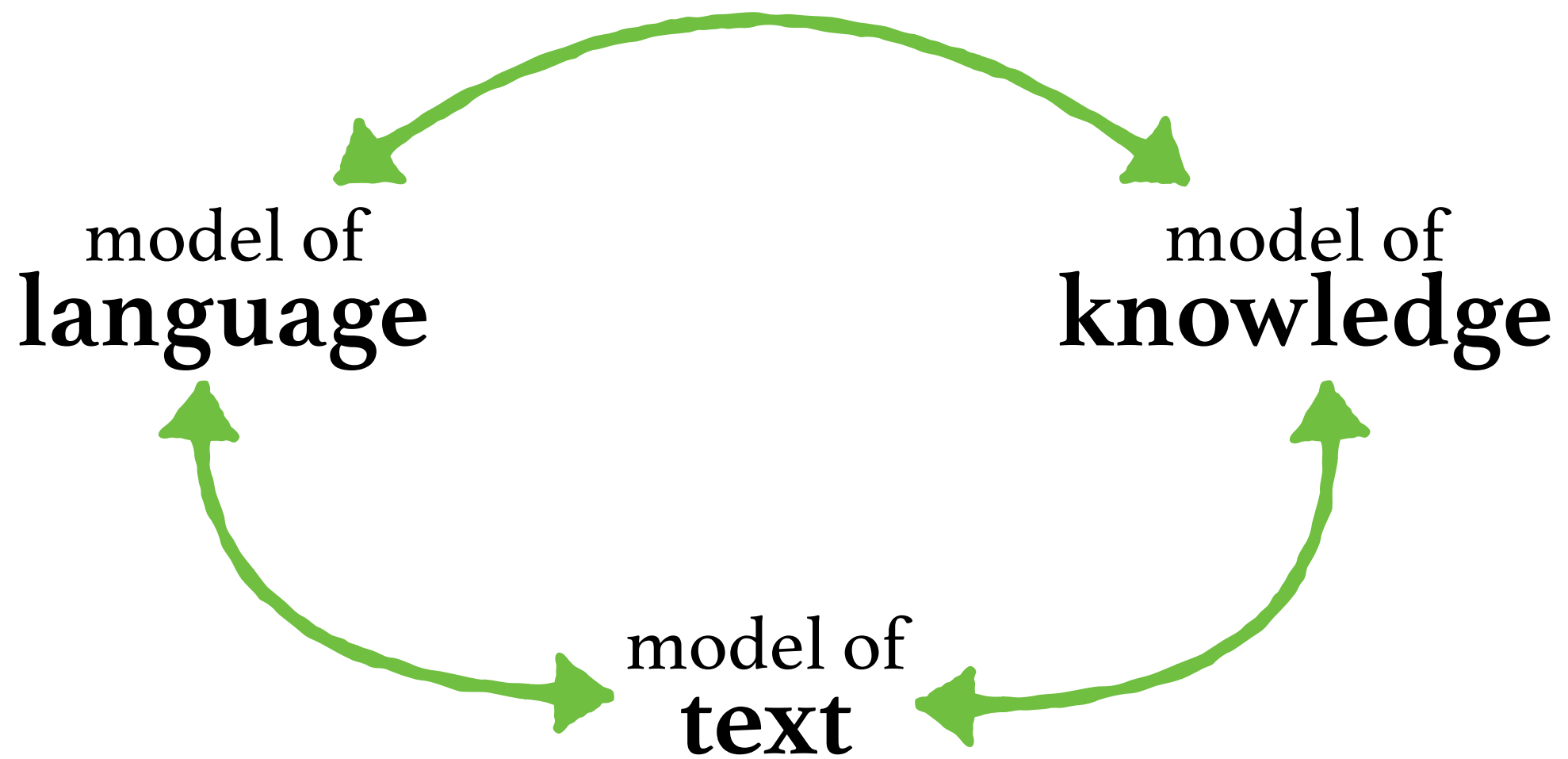
Γαλατία:Γαλατικός

Πόντος:Ποντικός

Στοῖκός

εἰρήνη:εἰρηνικός

ὄνος:ὄνικός



Adaptive Reading Environments

what's needed to understand an upcoming passage

what the student has already seen

what the student has inquired about

what is at an optimal recall interval

what the student is good or not so good at
understanding (based on explicit assessment
including meta-cognitive questions)

[Log in](#) [Sign up](#)



Perseus Digital Library Scaife Viewer

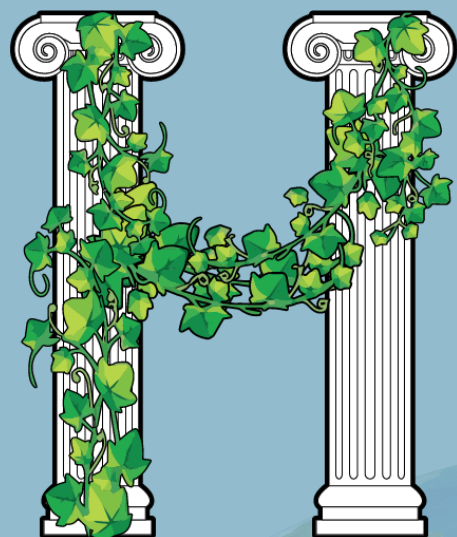
[Browse Library](#)

[Text Search](#)

or try the Iliad or Plato's Apology.

1,834 works in **2,320** editions and translations (1,178 in Greek and 612 in Latin)

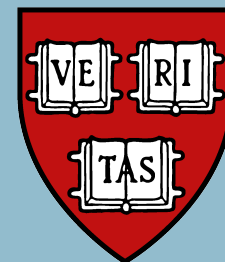
50.1 million words (20.1 million in Greek, 15.6 million in Latin)



EDERA

reading that fits

<http://hederaproject.org>



HARVARD
DEPARTMENT OF THE CLASSICS
ACADEMIC
TECHNOLOGY
for the FACULTY OF ARTS & SCIENCES

Teachers

"How hard will this be for my students to read?"

Lemmatized Texts » demo3

hic rex hostes vincere optat. milites suos vocat et hortatur: "animos monstrate! plus honoris et plus virtutis quam illi habetis. cogitate de uxoribus et **filiabus** vestros. nunc arma capite atque patriam defendite!" clamore sublato, milites arma sustulerunt. acie instructo, milites contra hostes celeriter cucurrerunt. quam acerrime pugnabant. tandem rex ducem hostium interfecit. duce interfecto et spe amissa, hostes ab eo loco fugerunt. rege laudato, milites gaudebant et domum redierunt. postea omnes vitam beatam agebant.

DCC Latin Core Vocabulary
The thousand most common words in Latin compiled by a team at Dickinson College led by Christopher Francese. See <http://dcc.dickinson.edu/vocab/core-vocabulary>

94.6%
Known
Select
The known words are highlighted.
Highlight Unknown

1944

form: filiae

2316

form: filiam

1346

lemma: filia

1347

lemma: filius

filia -ae f.;
filius -i m.

Life-Long Learners

"What would be an appropriate text to read next?"

Text	Language	Length	Lemmatized	Familiarity
Livy, Ad Urbe Condita Teacher • Learner	lat	55	<div></div>	44.19% 0.00% 0.00% 6.98% 11.63% 37.21%
01 Ecce Aeneas Teacher • Learner	lat	158	<div></div>	72.73% 0.00% 2.02% 3.03% 8.08% 43.43%
02 Amulius et Numitor Teacher • Learner	lat	174	<div></div>	71.13% 0.00% 2.06% 3.09% 7.22% 41.24%

Students

"How well do I know these words?"

Lemmatized Texts » Livy, Ad Urbe Condita

Proca deinde regnat. Is **Numitorem** atque **Amulium** procreat, **Numitori**, qui stirpis maximus erat, regnum vetustum **Silviae** gentis legat. Plus tamen vis potuit quam voluntas patris aut verecundia aetatis: pulso fratre **Amulius** regnat. Addit sceleri scelus: stirpem fratris virilem interemit, fratris filiae Reae **Silviae** per speciem honoris cum **Vestalem** eam legisset perpetua virginitate spem partus adimit.

55 Total Tokens
43 Unique Tokens

44.2% 0.0% 0.0% 7.0% 11.6% 37.2%

is ea id
Reveal Gloss

Show Familiarity

Summary

language acquisition needs lots of comprehensible input

we can sequence texts and vocabulary hand-in-hand to achieve this (and consider morphology, syntactic constructions, etc too)

this all relies on richly annotated texts (and other machine-actionable language description)

it can be done statically but the real power comes in adapting to students as they learn

jktauber.com