



Input Sparsity and Derivational Relationships in Latin and Spanish

Jordan Kodner

Stony Brook University

Productivity Workshop at DGfS 44, Feb 24, 2022

Outline

- **Classical Latin Past Participles and Derivatives**
- **Spanish Past Participles and Derivatives**
- **Language Acquisition and Productivity**
- **Explaining the Spanish System with Diachrony**



The situation in Latin

Classical Latin Principal Parts and Conjugations

- Traditionally classified into 4½ conjugations distinguished by 4 principal parts
- Conjugations correspond to theme vowels, principal parts to stems

Principal parts

1. present active indicative 1sg
2. present active infinitive
3. perfect active indicative 1sg
4. past participle (or supine)

Conj.	ThV	1st PP present stem	2nd PP	3rd PP perfect	4th PP pptc	Meaning
1st	<i>ā</i>	<i>amō</i>	<i>amāre</i>	<i>amāvī</i>	<i>amātus</i>	‘love’
2nd	<i>ē</i>	<i>moneō</i>	<i>monēre</i>	<i>monuī</i>	<i>monitus</i>	‘warn’
3rd	<i>e</i>	<i>legō</i>	<i>lēgere</i>	<i>lēgī</i>	<i>lēctus</i>	‘choose’
3rd -iō	<i>i</i>	<i>capiō</i>	<i>capere</i>	<i>cēpī</i>	<i>captus</i>	‘take’
4th	<i>ī</i>	<i>audiō</i>	<i>audīre</i>	<i>audīvī</i>	<i>audītus</i>	‘hear’

Complex Forms of the Past Participle

- Stems are not reliably derivable from one another

	Present	Perfect	PPtc	Meaning
<i>amō</i>	<i>amāre</i>	<i>amāvī</i>	<i>amātus</i>	‘love’
<i>sonō</i>	<i>sonāre</i>	<i>sonuī</i>	<i>sonitus</i>	‘sound’
<i>moneō</i>	<i>monēre</i>	<i>monuī</i>	<i>monitus</i>	‘warn’
<i>maneō</i>	<i>manēre</i>	<i>mānsī</i>	<i>mānsus</i>	‘stay’
<i>teneō</i>	<i>tenēre</i>	<i>tenuī</i>	<i>tentus</i>	‘hold’
<i>audiō</i>	<i>audire</i>	<i>audīvī</i>	<i>audītus</i>	‘hear’
<i>pellō</i>	<i>pellere</i>	<i>pepulī</i>	<i>pulsus</i>	‘push’
<i>capiō</i>	<i>capere</i>	<i>cēpī</i>	<i>captus</i>	‘take’
<i>ferō</i>	<i>ferre</i>	<i>tulī</i>	<i>lātus</i>	‘carry’

Complex Forms of the Past Participle

- Stems are not reliably derivable from one another

Verbs with similar forms for one stem may not have similar forms for the others

Present		Perfect	PPtc	Meaning
<i>amō</i>	<i>amāre</i>	<i>amāvī</i>	<i>amātus</i>	'love'
<i>sonō</i>	<i>sonāre</i>	<i>sonuī</i>	<i>sonitus</i>	'sound'
<i>moneō</i>	<i>monēre</i>	<i>monuī</i>	<i>monitus</i>	'warn'
<i>maneō</i>	<i>manēre</i>	<i>mānsī</i>	<i>mānsus</i>	'stay'
<i>teneō</i>	<i>tenēre</i>	<i>tenuī</i>	<i>tentus</i>	'hold'
<i>audiō</i>	<i>audire</i>	<i>audīvī</i>	<i>audītus</i>	'hear'
<i>pellō</i>	<i>pellere</i>	<i>pepulī</i>	<i>pulsus</i>	'push'
<i>capiō</i>	<i>capere</i>	<i>cēpī</i>	<i>captus</i>	'take'
<i>ferō</i>	<i>ferre</i>	<i>tulī</i>	<i>lātus</i>	'carry'

5 forms
7 forms
7 forms

Conjugations and PPTcs by Type Count

Data extracted from all the Old and Classical Latin from Perseus¹

- ~3.5 million tokens
- POS-tagged and lemmatized with modified CLTK²

Conjugation	# Verbs	Top freq	% Top	Next most	% Top two
1st	541	<i>-ātus</i> 528	97.6%	<i>-itus</i> 6	98.7%
2nd	65	<i>-itus</i> 25	38.5%	<i>-tus</i> 17	64.6%
3rd	215	<i>-tus</i> 80	37.2%	<i>-itus</i> 19	46.6%
4th	55	<i>-ītus</i> 34	61.8%	<i>-tus</i> 13	87.3%

¹ Smith et al (2020), ² <http://cltk.org/>

Conjugations and PPTcs by Type Count

Out of the most frequent verbs,

- 1st conjugation is largest and most homogeneous

Conjugation	# Verbs	Top freq	% Top	Next most	% Top two
1st	541	-ātus 528	97.6%	-itus 6	98.7%
2nd	65	-itus 25	38.5%	-tus 17	64.6%
3rd	215	-tus 80	37.2%	-itus 19	46.6%
4th	55	-ītus 34	61.8%	-tus 13	87.3%

Conjugations and PPTcs by Type Count

Out of the most frequent verbs,

- 1st conjugation is largest and most homogeneous
- 3rd conjugation is second largest and most heterogeneous

Conjugation	# Verbs	Top freq	% Top	Next most	% Top two
1st	541	<i>-ātus</i> 528	97.6%	<i>-itus</i> 6	98.7%
2nd	65	<i>-itus</i> 25	38.5%	<i>-tus</i> 17	64.6%
3rd	215	<i>-tus</i> 80	37.2%	<i>-itus</i> 19	46.6%
4th	55	<i>-ītus</i> 34	61.8%	<i>-tus</i> 13	87.3%

Conjugations and PPTcs by Type Count

Out of the most frequent verbs,

- 1st conjugation is largest and most homogeneous
- 3rd conjugation is second largest and most heterogeneous
- *-itus* and *-tus* are the most common pptcs outside the 1st conjugation

Conjugation	# Verbs	Top freq	% Top	Next most	% Top two
1st	541	<i>-ātus</i> 528	97.6%	<i>-itus</i> 6	98.7%
2nd	65	<i>-itus</i> 25	38.5%	<i>-tus</i> 17	64.6%
3rd	215	<i>-tus</i> 80	37.2%	<i>-itus</i> 19	46.6%
4th	55	<i>-ītus</i> 34	61.8%	<i>-tus</i> 13	87.3%

Conjugations and PPTcs by Type Count

Out of the most frequent verbs,

- 1st conjugation is largest and most homogeneous
- 3rd conjugation is second largest and most heterogeneous
- *-itus* and *-tus* are the most common pptcs outside the 1st conjugation

Conjugation	# Verbs	Top freq	% Top	Next most	% Top two
1st	541	<i>-ātus</i> 528	97.6%	<i>-itus</i> 6	98.7%
2nd	65			<i>-sus</i> 17	64.6%
3rd	215			<i>-tus</i> 19	46.6%
4th	55	<i>-ītus</i> 34	61.8%	<i>-tus</i> 13	87.3%

What counts as regular here?

The Classical Latin *t*-Deverbals

- Deverbals with suffixes containing *t* (or *s*)
- A wide range of syntactic categories and meanings

Type	Ending	Verb	Meaning	<i>t</i> -Deverbal	Meaning
Adverb	<i>-tim</i>	<i>stō</i>	‘stand’	<i>statim</i>	‘immediately’
Agent	<i>-tor</i>	<i>doceō</i>	‘teach’	<i>doctor</i>	‘teacher’
Event	<i>-tiō</i>	<i>agō</i>	‘do’	<i>actiō</i>	‘action’
Event	<i>-tus</i>	<i>sūmō</i>	‘spend’	<i>sumptus</i>	‘expenditure’
Fut Ptc	<i>-tūrus</i>	<i>morior</i>	‘die’	<i>moritūrus</i>	‘about to die’
Result	<i>-tūra</i>	<i>scribō</i>	‘write’	<i>scriptūra</i>	‘writing’

Derivation of the Classical *t*-Deverbals

- *t*-Deverbals appear to be constructed from the pptic stem
- They adopt whatever irregularities exist in the pptic, including suppletion

Priscian Algorithm

- Begin with pptic
- Delete case/number ending
- Add *t*-deverbal ending
- Done!

1st PP	2nd PP	3rd PP	4th PP	<i>t</i> -Deverbal
<i>amō</i>	<i>amāre</i>	<i>amāvī</i>	<i>amātus</i>	<i>amātor</i>
<i>habeō</i>	<i>habēre</i>	<i>habuī</i>	<i>habitus</i>	<i>habitor</i>
<i>agō</i>	<i>agere</i>	<i>ēgī</i>	<i>actus</i>	<i>actor</i>
<i>pellō</i>	<i>pellere</i>	<i>pepulī</i>	<i>pulsus</i>	<i>pulsor</i>
<i>sequor</i>	<i>sequī</i>	<i>secūtus est</i>	-	<i>secūtor</i>
<i>ferō</i>	<i>ferre</i>	<i>tulī</i>	<i>lātus</i>	<i>lātor</i>

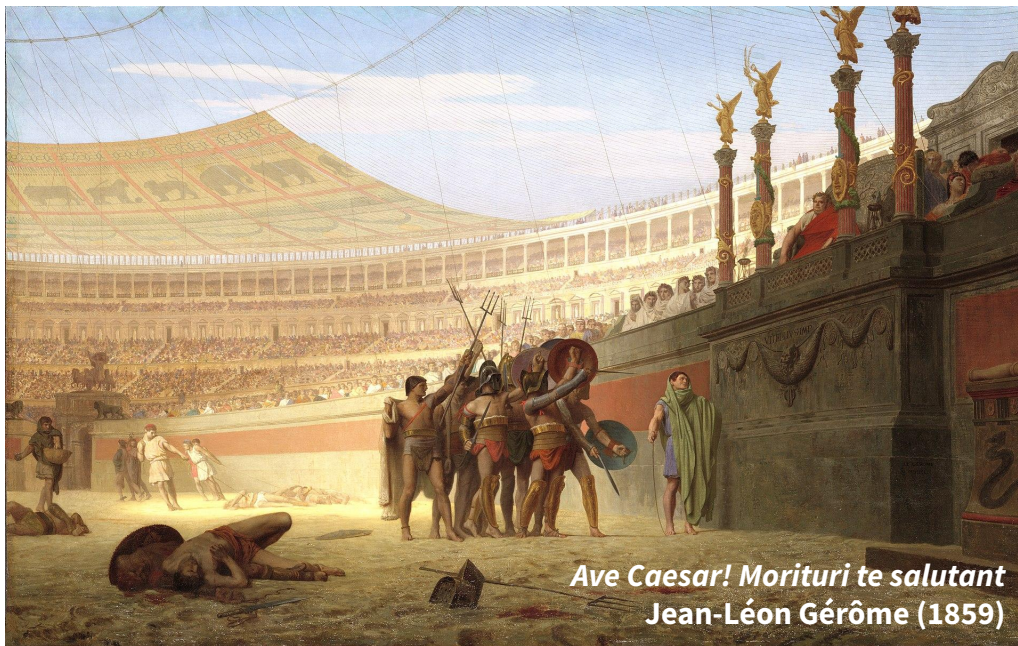
Derivation of the Classical *t*-Deverbals

- *t*-Deverbals appear to be constructed from the pptic stem
- They adopt whatever irregularities exist in the pptic, including suppletion

Exceptions are limited

- PPtc *mortuus* ‘dead’
FPtc *moritūrus* ‘about to die’
- PPtc *sonituus* ‘sounded’
FPtc *sonātūrus* ‘about to sound’
- PPtc *lautus* ‘washed’
Supine *lavātum*

The correspondence
is productive with
few exceptions



Practical Productivity

PPTc predicts *t*-deverbal and vice-versa

- PPTcs are far more common than any *t*-deverbal in the corpus
- In practice, inference pptc → *t*-deverbal is much more common

Category	#Freq ≥ 35	%Total	#Unique	% of Category	% of Unique
PPTc	1006	75.9%	817	81.2%	89.6%
Adverb	18	1.4%	8	44.4%	0.9%
Agent	72	5.4%	20	27.7%	2.2%
Event	178	13.4%	54	30.3%	5.9%
FPtc	52	3.9%	13	25.0%	1.5%
Total	1326		912	68.8%	

Practical Productivity

PPTc predicts *t*-deverbal and vice-versa

- PPTcs are far more common than any *t*-deverbal in the corpus
- In practice, inference pptc → *t*-deverbal is much more common

How many *t*-devs are at least as frequent as the 1000th most freq pptc?

Category	#Freq ≥ 35	%Total	#Unique	% of Category	% of Unique
PPTc	1006	75.9%	817	81.2%	89.6%
Adverb	18	1.4%	8	44.4%	0.9%
Agent	72	5.4%	20	27.7%	2.2%
Event	178	13.4%	54	30.3%	5.9%
FPtc	52	3.9%	13	25.0%	1.5%
Total	1326		912	68.8%	

Practical Productivity

PPTc predicts *t*-deverbal and vice-versa

- PPTcs are far more common than any *t*-deverbal in the corpus
- In practice, inference pptc → *t*-deverbal is much more common

How many *t*-devs are at least as frequent as the 1000th most freq pptc?

How many stems are attested only in a *t*-dev or only the pptc?

Category	#Freq ≥ 35	%Total	#Unique	% of Category	% of Unique
PPTc	1006	75.9%	817	81.2%	89.6%
Adverb	18	1.4%	8	44.4%	0.9%
Agent	72	5.4%	20	27.7%	2.2%
Event	178	13.4%	54	30.3%	5.9%
FPtc	52	3.9%	13	25.0%	1.5%
Total	1326		912	68.8%	



The situation in Spanish

Spanish Past Participles

- Three conjugations (-ar < Lat. -āre, -er < Lat. -ēre, -ir < Lat. -ere and -īre)
- Past participles are highly regular but not exceptionless

Conj.	Present	Preter.	PPtc	Latin	Meaning
-ar	<i>amar</i>	<i>amé</i>	<i>amado</i>	< <i>amāt-</i>	‘love’
-er	<i>vencer</i>	<i>vencí</i>	<i>vencido</i>	(<i>vict-</i>)	‘defeat’
-ir	<i>sentir</i>	<i>sentí</i>	<i>sentido</i>	(<i>sēns-</i>)	‘feel’
irreg	<i>hacer</i>	<i>hice</i>	<i>hecho</i>	< <i>fāct-</i>	‘make’
irreg	<i>ver</i>	<i>vi</i>	<i>visto</i>	< (<i>vīs-</i>)	‘see’
irreg	<i>escribir</i>	<i>escribí</i>	<i>escrito</i>	< <i>script-</i>	‘write’

Spanish Past Participles

- Three conjugations (-ar < Lat. -āre, -er < Lat. -ēre, -ir < Lat. -ere and -īre)
- Past participles are highly regular but not exceptionless

	Conj.	Present	Preter.	PPtc	Latin	Meaning
Inherited from Latin Reworked on basis of present	-ar	<i>amar</i>	<i>amé</i>	<i>amado</i>	< <i>amāt-</i>	‘love’
	-er	<i>vencer</i>	<i>vencí</i>	<i>vencido</i>	(<i>vict-</i>)	‘defeat’
	-ir	<i>sentir</i>	<i>sentí</i>	<i>sentido</i>	(<i>sēns-</i>)	‘feel’
Inherited from Latin	irreg	<i>hacer</i>	<i>hice</i>	<i>hecho</i>	< <i>fāct-</i>	‘make’
	irreg	<i>ver</i>	<i>vi</i>	<i>visto</i>	< (<i>vīs-</i>)	‘see’
	irreg	<i>escribir</i>	<i>escribí</i>	<i>escrito</i>	< <i>script-</i>	‘write’

Spanish Past Participles

- Three conjugations (*-ar* < Lat. *-āre*, *-er* < Lat. *-ēre*, *-ir* < Lat. *-ere* and *-īre*)
- Past participles are highly regular but not exceptionless

Conjugation	# PPtcs	Reg. PPtc		% Reg
<i>-ar</i>	373	<i>-ado</i>	373	100%
<i>-er</i>	70	<i>-ido</i>	58	82.9%
<i>-ir</i>	94	<i>-ido</i>	81	86.2%

Spanish *t*-Deverbals

- *t*-Deverbal agent nouns and event nouns survive from Latin
- But note *-ción* is itself borrowed form Latin (doublet with inherited *-zón*)

Type	Verb	<i>t</i> -Deverbal	Meaning	Latin <i>t</i> -Dev	Meaning
Agent	<i>amar</i>	<i>amador</i>	‘lover’	<i>amātor</i>	‘lover’
Agent	<i>vencer</i>	<i>vencedor</i>	‘conqueror’	<i>victor</i>	‘conqueror’
Agent	<i>batir</i>	<i>batidor</i>	‘whisk’	(none)	‘beat’
Event	<i>quemar</i>	<i>quemazón</i>	‘burning’	<i>cremātion-</i>	‘burning’
Event	<i>comer</i>	<i>comezón</i>	‘itching’	<i>comestion-</i>	‘eating’
Event	<i>partir</i>	<i>partición</i>	‘partition’	<i>partition-</i>	‘distribution’

Spanish *t*-Deverbals

- *t*-Deverbals correspond to the present rather than pptc if different
- Irregular *t*-deverbals are usually borrowed from Latin

	Verb	PPtc	<i>t</i> -Deverbal	Latin PPtc	Meaning
Reworked on basis of present	<i>hacer</i>	<i>hecho</i>	<i>hacedor</i>	<i>fāct-</i>	‘maker’
	<i>abrir</i>	<i>abierto</i>	<i>abridor</i>	<i>apert-</i>	‘opener’
Borrowed from Latin	<i>poner</i>	<i>puesto</i>	<i>posición</i>	<i>posit-</i>	‘position’
	<i>devolver</i>	<i>devuelto</i>	<i>devolución</i>	<i>-volūt-</i>	‘devolution’
	<i>leer</i>	<i>leído</i>	<i>lección</i>	<i>lēct-</i>	‘lesson’
	<i>conducir</i>	<i>conducido</i>	<i>conductor</i>	<i>con-dūct-</i>	‘driver’

Practical Productivity

PPTc predicts *t*-deverbal and vice-versa

- PPTcs are more common than any *t*-deverbal in the corpus
- In practice, inference pptc → *t*-deverbal is more common
- But this is **far less skewed than Latin**
- Many event nouns seem to be borrowed rather than synchronically derived

How many *t*-devs are at least as frequent as the 500th most freq pptc?

How many stems are attested only in a *t*-dev or only the pptc?

Category	#Freq ≥ 7	%Total	#Unique	% of Category	% of Unique
PPTc	540	54.7%	408	75.6%	61.0%
Agent	105	10.6%	52	49.5%	7.8%
Event	342	34.7%	209	61.1%	31.2%
Total	986		669	69.2%	

Practical Productivity

PPTc predicts *t*-deverbal and vice-versa

- PPTcs are more common than any *t*-deverbal in the corpus
- In practice, inference pptc → *t*-deverbal is more common
- But this is **far less skewed than Latin**
- Many event nouns seem to be borrowed rather than synchronically derived

Category	#Freq	%Total	#Unique	% of Category	% of Unique
ES PPTc	540	54.7%	408	75.6%	61.0%
LA PPTc	1006	75.9%	817	81.2%	89.6%

Interim Summary

Classical Latin

- **Complex relationship between pptc and other stems**
- ***t*-Devs correspond to pptcs regardless of pptc regularity**
- **PPTcs are much more frequent than all *t*-devs combined**

Modern Spanish

- **PPTcs almost always predictable from present stem**
- ***t*-Devs correspond with the present even if pptc is irregular**
- **PPTcs are more frequent than *t*-devs but not as skewed as Latin**

Productivity, Learning, and Change

Leveraging Child Language Acquisition

- **Determination of productive patterns is a central question in acquisition**
- **Exemplified by the English “Past Tense Debate”¹**
 - How are patterns and exceptions learned?
 - How are developmental trajectories explained?

¹ Rumelhart & McClelland 1986, Pinker & Prince 1988, Pinker 1994, Albright & Hayes 2006, Yang 2005, *and many more*

Leveraging Child Language Acquisition

- **Determination of productive patterns is a central question in acquisition**
- **Exemplified by the English “Past Tense Debate”¹**
 - How are patterns and exceptions learned?
 - How are developmental trajectories explained?

Broad agreement:

it isn't just token frequency (and derived measures)!²

→ **Quantitative corpus analysis alone won't cut it**

→ **Should work through the implications of some concrete learning mechanism**

¹ Rumelhart & McClelland 1986, Pinker & Prince 1988, Pinker 1994, Albright & Hayes 2006, Yang 2005, *and many more*

² Aronoff 1976, MacWhinney 1978, Bybee 1985, Baayen 1993, Elman 1998, Pierrehumbert 2003, Yang 2016

The Tolerance Principle

- An **evaluation metric**¹ over linguistic hypotheses
- Is derived from
 - an **Elsewhere Condition** for ‘rules’ and ‘exceptions’²
 - **frequency-rank correlated lexical access**³
 - Generally **Zipfian** input distributions
- Received **psychological backing from artificial language learning experiments**⁴

¹Chomsky 1955, 1965, Chomsky & Halle 1968, ²Anderson 1969, *inter alia*, ³Murray & Forster 2004, ⁴Schuler et al 2017, Emond & Shi 2020

The Tolerance Principle

- An **evaluation metric**¹ over linguistic hypotheses
- Is derived from
 - an **Elsewhere Condition** for ‘rules’ and ‘exceptions’²
 - **frequency-rank correlated lexical access**³
 - Generally **Zipfian** input distributions
- Received **psychological backing** from artificial language learning experiments⁴

Example Applications

- Is **-s** the default German noun pl? Under what conditions is **-(e)n** productive?
- Is vowel mutation as in **sing~sang** productive among similar verbs?

¹Chomsky 1955, 1965, Chomsky & Halle 1968, ²Anderson 1969, *inter alia*, ³Murray & Forster 2004, ⁴Schuler et al 2017, Emond & Shi 2020

The Tolerance Principle¹

Given a hypothesized generalization R operating over a class C , quantitatively define the number of exceptions below which the generalization is tenable

The Tolerance Principle¹

Given a hypothesized generalization R operating over a class C , quantitatively define the number of exceptions below which the generalization is tenable

N = number of **types** that should obey the generalization

e = number of **types** that **do not** obey the generalization

θ = max # of exceptions that can be tolerated

Exceptions are **tolerable** if

$$e < \theta$$

$$\theta = N / \ln N$$

Visualization of the Tolerance Principle

N = types it should apply to

e = types that are exceptions

θ = tolerance threshold



e falls in the range $[0, N]$ and may be less than or greater than θ

Visualization of the Tolerance Principle

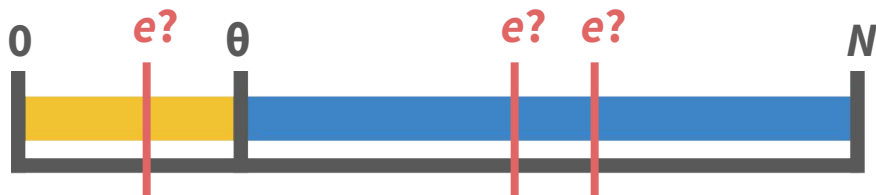
N = types it should apply to
 e = types that are exceptions
 θ = tolerance threshold



If e is below θ ,
acquire generalization

Visualization of the Tolerance Principle

N = types it should apply to
 e = types that are exceptions
 θ = tolerance threshold



If e is below θ ,
acquire generalization
Otherwise, do not generalize

Visualization of the Tolerance Principle

N = types it should apply to
 e = types that are exceptions
 θ = tolerance threshold

If e is below θ ,
acquire generalization

Otherwise, do not generalize

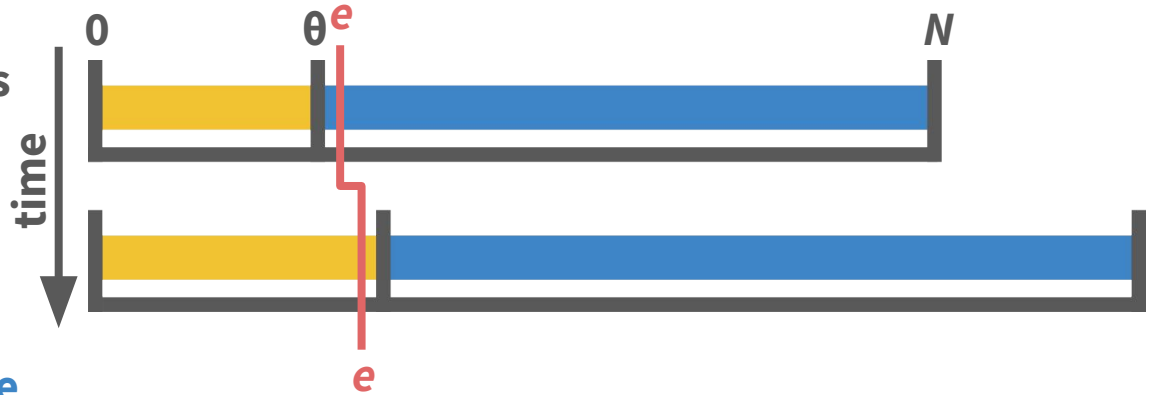


- N grows over an individual's development, θ grows more slowly

Visualization of the Tolerance Principle

N = types it should apply to
 e = types that are exceptions
 θ = tolerance threshold

If e is below θ ,
acquire generalization
Otherwise, do not generalize



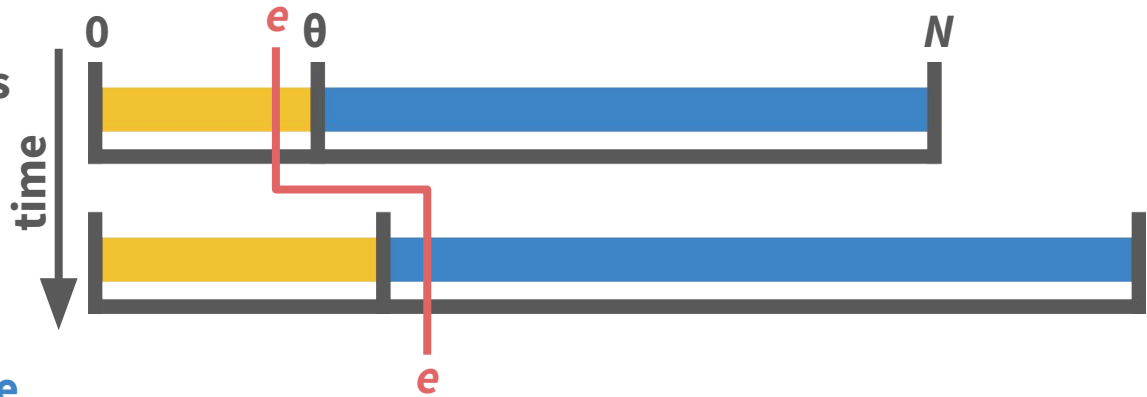
- N grows over an individual's development, θ grows more slowly
- If θ grows faster than e , a generalization may fall into productivity

Visualization of the Tolerance Principle

N = types it should apply to
 e = types that are exceptions
 θ = tolerance threshold

If e is below θ ,
acquire generalization

Otherwise, do not generalize



- N grows over an individual's development, θ grows more slowly
- If θ grows faster than e , a generalization may fall into productivity
- If e grows faster than θ , a generalization may fall out of productivity

Child Lexical Knowledge

- Learners' vocabularies grow over the course of development
- There is significant individual variation, but consistent trends
- Only on the order of 10^2 for English and German learners by around age 3
- Children have the foundations for language-specific grammars by this point

A roughly 1 per million frequency cutoff applied to the larger CHILDES corpora yields lexicons like these¹

Language	Estimated Vocab
English 2;10-3;0 ²	525-1,116
German 2;6 ³	$\mu = 429, \sigma > 100$

¹ Nagy & Anderson 1984, ² Hart & Risley 2003, ³ Szagun et al 2006

Applying the Tolerance Principle

Over likely generalizations

- Present stem → *t*-dev forms
- PPTc stem → *t*-dev forms

Theory independent interpretation

- Generalizations over surface phonotactics “rightmost vowel is /a:/”
- Or generalizations over morphemes “ThV is *-ā-*”

Example Calculation

Is *stem+āt-* the productive *t*-dev for verbs with Theme V *ā*?

Example Calculation

Is *stem+āt-* the productive *t*-dev for verbs with Theme V *ā*?

A typical child who knows $n=500$ verbs knows

- $N=221$ *ā* verbs
- $e=13$ *ā* verbs with non *-āt-* *t*-devs

Example Calculation

Is *stem+āt-* the productive *t*-dev for verbs with Theme V *ā*?

A typical child who knows $n=500$ verbs knows

- $N=221$ *ā* verbs
- $e=13$ *ā* verbs with non *-āt-* *t*-devs

Exceptions are tolerable if

$$e < \theta$$

$$\theta = N / \ln N$$

Example Calculation

Is *stem+āt-* the productive *t*-dev for verbs with Theme V *ā*?

A typical child who knows $n=500$ verbs knows

- $N=221$ *ā* verbs
- $e=13$ *ā* verbs with non *-āt-* *t*-devs
- $\theta=40.94$ tolerance threshold

Exceptions are tolerable if

$$13 < 40.9$$

$$\theta = N / \ln N$$

-āt- is productive for *ā* verbs at $n=500$

Summary results for Past Participles¹

If derivations are only possible from the present,

- Productive pptic derivation for 1st (-ātus), 3rd-*iō* (-tus)
- Marginal for *faveō*-type (-autus/-ōtus) and *solvō*-type (-ūtus)
- **No productive pptic derivation for 2nd, 3rd-*ō*, 4th**
- **No broadly productive -ītus or -tus**

If derivations is possible from the perfect,

- The above + productive deriv for -īvī (most of 4th; -ītus), -ēvī (-ētus), -Csī (-tus)
- Solidly productive -ūtus for *solvō*-types
- **No broadly productive pptic derivation for -uī-perfect verbs**
- **Still no broadly productive -ītus or -tus**

¹ Kodner (to appear)

The Past Participle ~ *t*-Deverbal Correspondence

Diachrony - It's mostly an accident

- The pptc and *t*-devs are etymologically related < PIE nominalizer **-to-*
- Same sound changes → same forms, eg *vīsus* ~ *vīsiō* < **wid-t-os*, **wid-t-iō-n-*
- **But not all forms are the result of sound change**, eg *offerō* ~ *oblātus* ~ *oblātio*

The Past Participle ~ *t*-Deverbal Correspondence

Diachrony - It's mostly an accident

- The pptic and *t*-devs are etymologically related < PIE nominalizer **-to-*
- Same sound changes → same forms, eg *vīsus* ~ *vīsiō* < **wid-t-os*, **wid-t-iō-n-*

Learning - Learning maintains the correspondence

- The form of most *t*-devs needs to be inferred - **sparsity problem**
- Most attested *t*-devs also have a corresponding attested pptic
- “**Make the *t*-dev be like the pptic**” works better than other hypotheses

Learning *t*-Deverbal Forms

Possible surface generalizations

1. Base the *t*-deverbals on the present or perfect
2. Base the *t*-deverbals on the pptc
3. Base the pptc on the *t*-dev

Methodology

- Test the generalizations on the Perseus corpus
- Using the Tolerance Principle

1. Base the *t*-Dev on the Present or Perfect

Works for a few classes, especially *-ā-* and *-ī-* stems

- Correspondence holds trivially for *-ā-* and *-ī-* stems
t-dev/pptc is thematic *-āt-* and *-īt-*
- Actually a majority of verbs!

But it doesn't work overall

- Too many exceptions for a learner to acquire

1. Base the *t*-Dev on the Present or Perfect

Works for a few classes, especially *-ā-* and *-ī-*stems

But it doesn't work overall

- Too many exceptions for a learner to acquire

Blue-Green

Productive

Red

Unproductive

Gold

within 1

White

Not evaluated

Theme Vowel	PPtc	Example	At <i>n</i> =100?	At 500?	At 1,000?
<i>ā</i> (1st)	<i>-ātus</i>	<i>vocāre ~ vocāt-</i>	<i>N</i> =17(<i>e</i> =1)	221 (11)	541 (14)
<i>ē</i> (2nd)	<i>-ītus</i>	<i>habēre ~ habit-</i>	16 (9)	55 (35)	65 (40)
<i>ē</i> (2nd)	<i>-tus</i>	<i>docēre ~ doct-</i>	16 (14)	55 (42)	65 (48)
<i>e</i> (3rd non- <i>iō</i>)	<i>-ītus</i>	<i>reddere~reddit-</i>	47 (46)	147 (136)	201(185)
<i>e</i> (3rd non- <i>iō</i>)	<i>-tus</i>	<i>scribere ~ script-</i>	47 (32)	147 (105)	201(143)
<i>i</i> (3rd - <i>iō</i>)	<i>-tus</i>	<i>capiō ~ captus</i>	9 (1)	12 (2)	14 (3)
<i>ī</i> (4th)	<i>-ītus</i>	<i>audire ~ audit-</i>	5 (3)	27 (9)	55 (21)
<i>ī</i> (4th)	<i>-tus</i>	<i>venire ~ vent-</i>	5 (2)	27 (20)	55 (42)

Individual Development



2. Base the *t*-Dev on the Past Participle

The correspondence overwhelmingly holds

- There are very few exceptions
- These tend to be high frequency → can be memorized

Some exceptions¹

- *mortuus* ‘dead’ but *moritūrus* ‘about to die’
- *sonitus* ‘sounded’ but *sonātūrus* ‘about to sound’

¹ Laurent 2003, pp. 18-19

3. Base the Past Participle on the *t*-Dev

In practice, inference has to go pptc → *t*-deverbal

- PPTcs are far more common than any *t*-deverbal in the corpus
- In practice, inference pptc → *t*-deverbal is much more common

3. Base the Past Participle on the *t*-Dev

In practice, inference has to go pptc → *t*-deverbal

- PPTcs are far more common than any *t*-deverbal in the corpus
- In practice, inference pptc → *t*-deverbal is much more common

How many *t*-devs are at least as frequent as the 1000th most freq pptc?

How many stems are attested only in a *t*-dev or only the pptc?

Category	#Freq ≥ 35	%Total	#Unique	% of Category	% of Unique
PPTc	1006	75.9%	817	81.2%	89.6%
Adverb	18	1.4%	8	44.4%	0.9%
Agent	72	5.4%	20	27.7%	2.2%
Event	178	13.4%	54	30.3%	5.9%
FPtc	52	3.9%	13	25.0%	1.5%
Total	1326		912	68.8%	

Correspondences in Spanish

PPTcs productively and transparently built on the present

- Much simpler than Latin
- Very few exceptions
- Conflation of *-er* and *-ir*

Conjugation	PPTc	Example	At <i>n</i> =500?
<i>-ar</i>	<i>-ado</i>	<i>amar ~ amado</i>	373 (0)
<i>-er</i>	<i>-ido</i>	<i>saber ~ sabido</i>	70 (12)
<i>-ir</i>	<i>-ido</i>	<i>seguir ~ seguido</i>	94 (13)

Agents almost always correspond with their pptcs

- But also agree with the present
- Agents *-er* shows *-e-* theme vowel *-edor*
- Agree with present over PPTc (eg *hacer, hecho, hacedor*)

Interim Summary

Classical Latin

- Complex relationship between pptc and other stems
- *t*-Devs correspond to pptcs regardless of pptc regularity
- PPTcs are much more frequent than all *t*-devs combined

Modern Spanish

- PPTcs almost always predictable from present stem
- *t*-Devs correspond with the present even if pptc is irregular
- PPTcs are more frequent than *t*-devs but less extreme than Latin
- PPTc→*t*-dev inference less important
- Ambiguous base for *t*-deverbal

Interim Summary

Classical Latin

- Complex relationship between pptc and other stems
- *t*-Devs correspond to pptcs regardless of pptc regularity
- PPTcs are much more frequent than all *t*-devs combined

Modern Spanish

- PPTcs almost always predictable from present stem
- *t*-Devs correspond with the present even if pptc is irregular
- PPTcs are more frequent than *t*-devs but less extreme than Latin
- PPTc→*t*-dev inference less important
- Ambiguous base for *t*-deverbal

How did the system realign from Latin to Spanish?

By type count, many Latin *t*-deverbals had an ambiguous base as well

Bridging Latin and Spanish

Remember how some Latin *t*-devs == pptcs == present stems?

- 1st conjugation is overwhelmingly regular Pres *-ā-* ~ PPtc *-āt-* ~ *t*-Dev *-āt-*
- Majority of 4th conj is too Pres *-ī-* ~ PPtc *-īt-* ~ *t*-Dev *-īt-*
- These support alternative Present ~ *t*-Dev analysis

Bridging Latin and Spanish

Remember how some Latin *t*-devs == pptcs == present stems?

- 1st conjugation is overwhelmingly regular Pres *-ā-* ~ PPtc *-āt-* ~ *t*-Dev *-āt-*
- Majority of 4th conj is too Pres *-ī-* ~ PPtc *-īt-* ~ *t*-Dev *-īt-*
- These support alternative Present ~ *t*-Dev analysis

The 1st and 4th conjugations grew in Late Latin

- Tendency to coin new intensive, iterative, etc, verbs in *-tāre*, *-titāre*, etc¹
 Inflect as “regular” first conjugation verbs **-atu* verbs build on present stems
 Replaced “irregular” 3rd conjugation verbs (eg *cantō*, *cantātus* replacing *canō*, *cantus*)²

Bridging Latin and Spanish

Remember how some Latin *t*-devs == pptcs == present stems?

- 1st conjugation is overwhelmingly regular Pres *-ā-* ~ PPtc *-āt-* ~ *t*-Dev *-āt-*
- Majority of 4th conj is too Pres *-ī-* ~ PPtc *-īt-* ~ *t*-Dev *-īt-*
- These support alternative Present ~ *t*-Dev analysis

The 1st and 4th conjugations grew in Late Latin

- Tendency to coin new intensive, iterative, etc, verbs in *-tāre*, *-titāre*, etc¹
 Inflect as “regular” first conjugation verbs **-atu* verbs build on present stems
 Replaced “irregular” 3rd conjugation verbs (eg *cantō*, *cantātus* replacing *canō*, *cantus*)²
- Distinction between 2nd, 3rd, and 4th conjugation collapsed variably into two categories *-er* < **-ere* < *-ēre* and *-ir* < **-ire* < *-īre*³
 These got “regular” pptcs build on present stems

Bridging Latin and Spanish

Remember how some Latin *t*-devs == pptcs == present stems?

- 1st conjugation is overwhelmingly regular Pres *-ā-* ~ PPtc *-āt-* ~ *t*-Dev *-āt-*
- Majority of 4th conj is too Pres *-ī-* ~ PPtc *-īt-* ~ *t*-Dev *-īt-*
- These support alternative Present ~ *t*-Dev analysis

The 1st and 4th conjugations grew in Late Latin

- Tendency to coin new intensive, iterative, etc, verbs in *-tāre*, *-titāre*, etc¹
 Inflect as “regular” first conjugation verbs **-atu* verbs build on present stems
 Replaced “irregular” 3rd conjugation verbs (eg *cantō*, *cantātus* replacing *canō*, *cantus*)²
- Distinction between 2nd, 3rd, and 4th conjugation collapsed variably into two categories *-er* < **-ere* < *-ēre* and *-ir* < **-ire* < *-īre*³
 These got “regular” pptcs build on present stems
- Spanish irregular pptcs are overwhelmingly high frequency and mostly inherited⁴ - what we expect from analogical leveling

Conclusions

Productivity in the *t*-deverbals over time

- **Derived from the past participle in Latin but present in Spanish**
- **Most Latin *t*-devs must be inferred from pptc**
But Spanish *t*-devs are more likely to be attested w/o the verb's pptc
- **Change in the past participles over time**
Largely unpredictable in Latin → Highly regular in Spanish
- **Modeling with the Tolerance Principle is consistent with this finding**

The End

Thank you

Four Features of Native Language Acquisition

1. **All children receive unique input yet exhibit gross developmental uniformity¹**
2. The type frequency of a pattern is crucial for acquisition of generalizations, as opposed to token frequency or attestation of initial items²
3. Token frequencies correlate with relative order of acquisition³
4. Early learner vocabularies are small⁴

¹ Labov 1972, ² Aronoff 1976, MacWhinney 1978, Bybee 1985, Baayen 1993, Elman 1998, Pierrehumbert 2003, Yang 2016, ³ Goodman 2008,

⁴ Hart & Risley 1995, 2003, Szagun et al. 2006

Four Features of Native Language Acquisition

1. All children receive unique input yet exhibit gross developmental uniformity¹
2. The type frequency of a pattern is crucial for acquisition of generalizations, as opposed to token frequency or attestation of initial items²
3. Token frequencies correlate with relative order of acquisition³
4. Early learner vocabularies are small⁴

¹ Labov 1972, ² Aronoff 1976, MacWhinney 1978, Bybee 1985, Baayen 1993, Elman 1998, Pierrehumbert 2003, Yang 2016, ³ Goodman 2008,

⁴ Hart & Risley 1995, 2003, Szagun et al. 2006

Four Features of Native Language Acquisition

1. All children receive unique input yet exhibit gross developmental uniformity¹
2. The type frequency of a pattern is crucial for acquisition of generalizations, as opposed to token frequency or attestation of initial items²
3. **Token frequencies correlate with relative order of acquisition³**
4. Early learner vocabularies are small⁴

¹ Labov 1972, ² Aronoff 1976, MacWhinney 1978, Bybee 1985, Baayen 1993, Elman 1998, Pierrehumbert 2003, Yang 2016, ³ Goodman 2008,

⁴ Hart & Risley 1995, 2003, Szagun et al. 2006

Four Features of Native Language Acquisition

1. All children receive unique input yet exhibit gross developmental uniformity¹
2. The type frequency of a pattern is crucial for acquisition of generalizations, as opposed to token frequency or attestation of initial items²
3. Token frequencies correlate with relative order of acquisition³
4. **Early learner vocabularies are small⁴**

¹ Labov 1972, ² Aronoff 1976, MacWhinney 1978, Bybee 1985, Baayen 1993, Elman 1998, Pierrehumbert 2003, Yang 2016, ³ Goodman 2008,

⁴ Hart & Risley 1995, 2003, Szagun et al. 2006

Four Features of Native Language Acquisition

1. All children receive unique input yet exhibit gross developmental uniformity¹
2. The type frequency of a pattern is crucial for acquisition of generalizations, as opposed to token frequency or attestation of initial items²
3. Token frequencies correlate with relative order of acquisition³
4. Early learner vocabularies are small⁴

As a result,

- Applying a frequency cutoff to lemmas in CDS approximates a “typical” child
- Insight taken by type frequency-based models of acquisition⁵

¹Labov 1972, ²Aronoff 1976, MacWhinney 1978, Bybee 1985, Baayen 1993, Elman 1998, Pierrehumbert 2003, Yang 2016, ³Goodman 2008,

⁴Hart & Risley 1995, 2003, Szagun et al. 2006, ⁵Nagy & Anderson 1984, Yang 2016

Acquisition in the Past

- Children in the past must have acquired language in the same way that modern children do - this is straightforward **uniformitarianism**¹
- We can reason about acquisition in the past in the same way we do now

Can non-CDS be substituted for CDS to study the relevant problem?

¹Labov 1972 as applied to linguistics, Walkden 2019

Acquisition in the Past

- Children in the past must have acquired language in the same way that modern children do - this is straightforward **uniformitarianism**¹
- We can reason about acquisition in the past in the same way we do now

Can non-CDS be substituted for CDS to study the relevant problem?

Yes, for the purposes of lexical acquisition²

¹Labov 1972 as applied to linguistics, Walkden 2019, ²Kodner 2019

Data Set

Perseus Corpus

- Scraped all Old and Classical Latin texts from website HTML
 - 3rd BC - AD 2nd inclusive
 - ~3.5mil tokens
- More than available by download - **undocumented “feature” :-**

Largest plain text OL/CL corpus?

Data Set

Perseus Corpus

- **Scraped all Old and Classical Latin texts from website HTML**
 - 3rd BC - AD 2nd inclusive
 - ~3.5mil tokens
- **More than available by download**

Post-Processing

- **POS-tagged and lemmatized with modified CLTK library**
 - 1,292 unique verb lemmas when derivational prefixes removed
- **Scraped Latin Wiktionary verbs to match lemmas to principal parts**

Data Set

Perseus Corpus

- **Scraped all Old and Classical Latin texts from website HTML**
 - 3rd BC - AD 2nd inclusive
 - ~3.5mil tokens
- **More than available by download**

Post-Processing

- **POS-tagged and lemmatized with modified CLTK library**
 - 1,292 unique verb lemmas when derivational prefixes removed
- **Scraped Latin Wiktionary verbs to match lemmas to principal parts**
- **Manually compared ~100 principal parts to Oxford Latin Dictionary**

Latin Wiktionary is surprisingly accurate!

Productive Present → PPtc by Theme Vowel

Theme Vowel	PPtc	Example	At n=100?	At 500?	At 1,000?
<i>ā</i> (1st)	<i>-ātus</i>	<i>vocāre ~ vocātus</i>	YES	YES	YES
<i>ē</i> (2nd)	<i>-ītus</i>	<i>habēre ~ habitus</i>	no	no	no
<i>ē</i> (2nd)	<i>-tus</i>	<i>docēre ~ doctus</i>	no	no	no
<i>e</i> (3rd non- <i>iō</i>)	<i>-ītus</i>	<i>reddere ~ redditus</i>	no	no	no
<i>e</i> (3rd non- <i>iō</i>)	<i>-tus</i>	<i>scribere ~ scriptus</i>	no	no	no
<i>i</i> (3rd - <i>iō</i>)	<i>-tus</i>	<i>capiō ~ captus</i>	YES	YES	YES
<i>e</i> or <i>i</i> (all 3rd)	<i>-ītus</i>	" ~ "	no	no	no
<i>e</i> or <i>i</i> (all 3rd)	<i>-tus</i>	" ~ "	no	no	no
<i>ī</i> (4th)	<i>-ītus</i>	<i>audīre ~ audītus</i>	YES	marginal*	no
<i>ī</i> (4th)	<i>-tus</i>	<i>venīre ~ ventus</i>	YES	no	no

Individual Development



* within 1 of threshold

Productive Present → PPTc more Narrowly

Present	PPTc	Example	At $n=100$?	At 500?	At 1,000?
-[a, o]veō	-[au, ō]tus	<i>faveō ~ faustus</i>	-	YES	YES
-[Velar]eō	-tus	<i>doceō ~ doctus</i>	-	no	no
-[not Velar]eō	-itus	<i>debeō ~ debitus</i>	marginal*	no	no
-[not Velar]eō	-tus	<i>teneō ~ tentus</i>	no	no	no
-vere	-ūtus	<i>solvere ~ solūtus</i>	YES	marginal*	marginal*
-[ll, rr]ere	-[l,r]sus	<i>currō ~ cursus</i>	-	marginal*	no
other 3rd	-ītus	<i>reddere ~ redditus</i>	no	no	no
other 3rd	-tus	<i>scribere ~ scriptus</i>	no	no	no

Individual Development



* within 1 of threshold

Productive Perfect → PPtC

Perfect	PPtc	Example	At n=100?	At 500?	At 1,000?
-āv-	-ātus	<i>amāvī ~ amātus</i>	YES	YES	YES
-īv-	-ītus	<i>dormīvī ~ dormītus</i>	YES	YES	YES
-ēv-	-ētus	<i>flēvī ~ flētus</i>	YES	YES	marginal*
-u-	-ītus	<i>valuī ~ valitus</i>	no	no	no
-u-	-tus	<i>tenuī ~ tentus</i>	no	no	no
-[Velar]u-	-tus	<i>liquī ~ lictus</i>	-	no	no
-[not Velar]u-	-ītus	<i>dēbuī ~ dēbitus</i>	no	no	no
-[not Velar]u-	-tus	<i>peruī ~ pertus</i>	no	no	no
-s-	-tus	<i>scripsī ~ scriptus</i>	no	no	no
-Cs-	-tus	<i>iūnxī ~ iūunctus</i>	YES	YES	YES
bare or stem change	-ītus	<i>lēgī ~ lēctus</i>	no	no	no

* within 1 of threshold

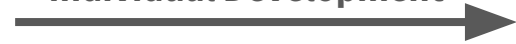
Individual Development



Productive Perfect + Present → PPtc

Perfect	PPtc	Example	At $n=100$?	At 500?	At 1,000?
-vere + -u-	-ūtus	<i>volvere ~ voluī ~ volūtus</i>	YES	YES	YES

Individual Development



- Only makes a difference for once class, but it is ***-utu**
- Only an option when a learner happens to know both stems

The System from Latin to Proto-Romance

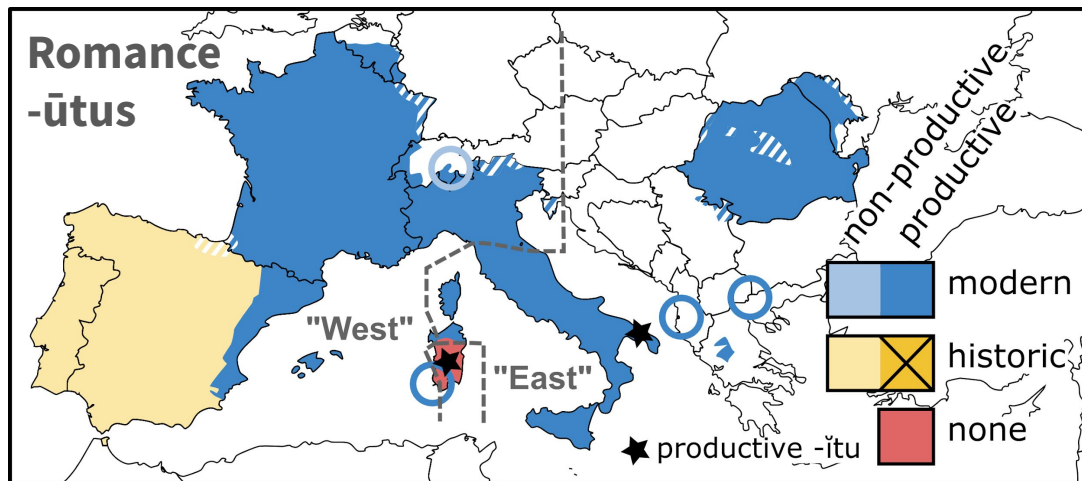
Varied across the Latin-speaking world, but in general...

- Novel verbs tended to have regular pptcs¹
- “Regular” **-atu*, **-itu*, **-utu* < *-ātus*, *-ītus* (not *-ītus*), *-ūtus* expanded at the expense of *-itus*, *-tus*, and others²
- The rise of **-utu* is mysterious given that it is rare in CL
- Perfects (→ preterites) were often regularized, often in **-ui* < *-uī*³

¹ Laurent 2003, ² *ibid.*, ³ *ibid.*

Reflexes of *-ūtus* and *-ītus* in Attested Romance¹

- Reflexives of *-ūtus* constitute the default for at least some class in most Romance languages
 - They are present but apparently non-productive in **Surselvan** (Rhaeto-Romance; Switzerland)
- Reflexes are attested in **Old Spanish** and **Portuguese** but have been lost
 - Their only reflexes are in adjectives eg, *agudo*, *menudo*
- *-ītus* remains productive in **Apulian** and **Sardinian**
 - /i/ merged with /i:/ in **Sardinian**, causing *-ītus* to fall together with *-ītus*



¹ data compiled from Laurent 2003

Diachronic Implications

Developments in Late Latin

- Three productive LL pptcs: **-atu* < *-ātus*, **-itu* < *-ītus*, **-utu* < *-ūtus*
- *-ītus* and *-tus* were unproductive in CL and reduced to irregulars
- *-ūtus* was productive for a small class
- But the only productive option for *-uī* perfects!
- It spread first among *-uī* perfects
- **No competition, “a big fish in a small pond”**

Implications

Listing and Rules

- **An externally motivated model guides theoretical analysis**
- **Predicts much more listing than a linguist relying on intuitions might**

The relationship between stems

- **If pptcs are derived from perfects**
 - **More can be derived by rule**
 - **Accounts for diachronic leveling of the perfect and pptc**
- **To do so, either perfect stems exist as representational objects or multiple step root → perfect “stem” → pptc derivations are required**

How are past participles derived?

- Are regular pptcs influenced by the present or perfect, or all memorized?
- Diachronic evidence for both
 - present → pptc: nasal infix spread
 - perfect → pptc: perfect analogies

The Nasal Infix

- Inherited from PIE, inserted into present stems
- Some continue to work like this in Latin¹
- But some have analogized to the perfect and pptc
- Only evidence for present → pptc derivation if absent in the perfect
 - At most two examples of this...
 - Otherwise, can present → perfect → pptc

Type	Present	Perfect	PPtc
Inherited	<i>fundō</i>	<i>fūdī</i>	<i>fūsus</i>
Pres, Perf	<i>fingō</i>	<i>fīnxī</i>	<i>fictus</i> ²
All	<i>iungō</i>	<i>iunxī</i>	<i>iūnctus</i>
Pres, PPtc	<i>pungō</i> <i>tundō</i>	<i>pupugī</i> <i>tutudī</i>	<i>pūnctus</i> <i>tū(n)sus</i>

¹ Poultney 1937, ² but Italian *finto*

Perfect Analogies

- Some pptcs have clearly been reworked on the basis of the perfect¹

cernō *crēvī* *crētus* (expected *certus* retained as adj)

sternō *strāvī* *strātus*

? *sonāre* *sonuī* *sonitus*

- Continues into Late Latin: eg **-utu* pptcs typically correspond to **-ui* perfects

¹ Table from Laurent 2003, p. 22

The System from Proto-Romance to Romance

Spanish, for example, shows the most regularization¹

- Regularization continued
 - *-ado*, *-ido*, and *-udo* existed in Old Spanish
 - Only *-ado*, *-ido* remain productive
- A handful of irregular pptcs remain, many relegated to adjectival meaning
 - *hecho*, *puesto*, *suelto*, *visto*, *vuelto*, etc, not all inherited
 - *teñir~teñido* ‘dyed’ but adj *tinto* ‘dyed red’ < *tinctus*, etc
 - OS had more eg *querer~quisto*, *prender~preso* < *prehensus*

¹ Laurent 2003 ch. 4.7

Past Participle Gaps and Meanings

- Past participles are typically passive
- But not all verbs have past participles¹
 - Sometimes due to semantics (eg, statives have no pptcs)
 - Sometimes they're more properly paradigmatic gaps

eg *bibō*, but *pōtus* not **bibitus*, *feriō*, but *percussus* not **ferītus*

- Some pptcs are active rather than passive²
 - Expected for deponents
 - But applies to some non-deponents as well

eg *locūtus* (deponent) 'having spoken,' *iūrātus* 'having sworn'

^{1,2} Laurent 2003, ² Embick 2000

Cross-Language Lexical Comparisons

- Compared lexical composition of modern CDS and historical corpora
- Calculated number of verb types across corpora with similar meanings

For corpus-derived lexicons A and B

where A and B are unordered sets,

$$\textit{similarity} = |A \cap B| / \min(|A|, |B|)$$

Cross-Language Corpora

- **English CDS** - verb lemmas in CHILDES Brown (and Brent for comparison)
- **Spanish CDS** - verb lemmas in combined CHILDES FernAguado, Hess, OreaPine, Remedi, Romero, SerraSole
- **Classical Latin** - verb lemmas in all Perseus online 3rd BC - 2nd AD (inclusive)

Corpus	Freq Cutoff	Lexicon size (<i>n</i>)
English CDS Brown	< 17	260
English CDS Brent	< 17	257
Spanish CDS	< 11	263
Latin	< 666	260

¹ Credit to Don Ringe for extracting them

Cross-Language Comparisons

- **Baselines: English-English (within-language) English-Spanish (cross-language)**
- **English-English unsurprisingly has the highest overlap**
- **Latin comparisons fall in between English-Spanish and English-English**

Latin Perseus contains the same kind of high frequency verbs that CDS does

Comparison	% Overlap
English - EN Brent	81.71%
English - Spanish	73.07%
English - Latin	75.77%
Spanish - Latin	78.62%

Paradigm Saturation

- **Paradigm Saturation¹** - the proportion of a verb's possible inflected forms which are actually attested in a corpus
- A measure of data sparsity
- Mean saturations tend to be low
- Obeys Zipfian distribution

¹ Chan 2008

Paradigm Saturation Data

- All POS-tagged, lemmatized, morpho feature annotated
- **CDS** - English (Brown), Spanish and German (CDS Leo¹)
- **Modern** - UD² English, Finnish, German, Spanish, Turkish
- **Historical** - UD Gothic, Latin
- Order 10⁵ verb tokens
- **CDS token/type ratios are on the order of 10x higher**

Corpus	Lang	# V Tokens	# V Types	Ratio
CDS	English	94,768	916	103.46
CDS	Spanish	96,686	879	110.00
CDS	German	81,351	641	126.91
Modern	English	53,796	3,225	16.67
Modern	Spanish	85,861	5,019	17.11
Modern	German	21,835	2,826	7.73
Modern	Finnish	63,891	3,476	18.38
Modern	Turkish	12,064	968	12.46
Historic	Gothic	12,749	1,172	10.88
Historic	Latin	99,066	2,2833	34.97

¹Behrens 2006, ²Nivre et al 2018

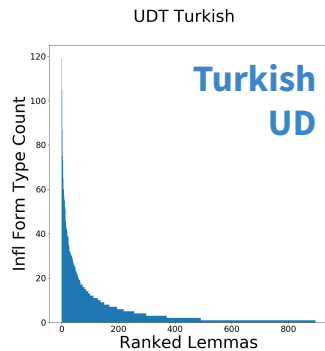
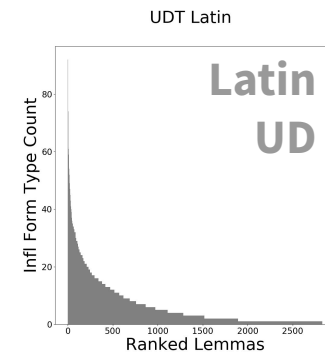
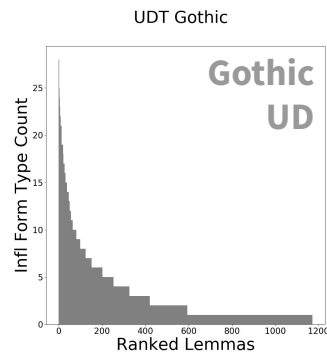
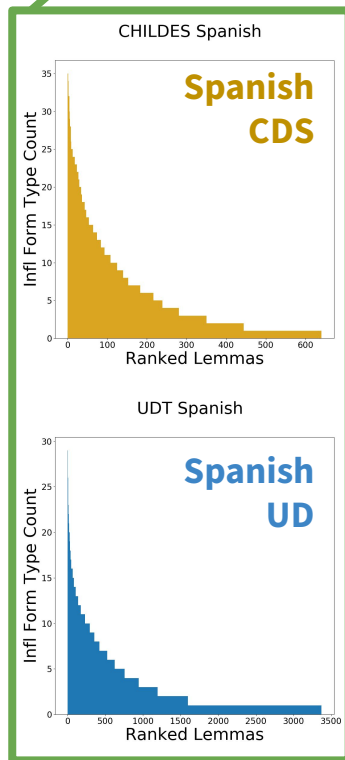
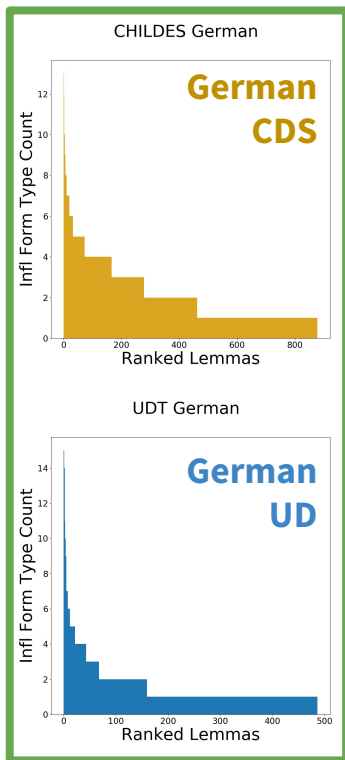
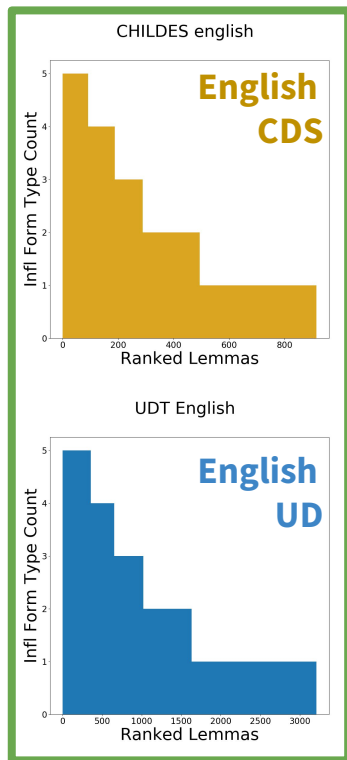
Paradigm Saturations

- **CDS saturations only slightly higher than modern equivs**
- **Despite difference in token/type ratios**
- **Historical corpora similar to modern ones**
- **Saturation appears related to paradigm size if anything**

Corpus	Lang	Paradigm	Max Sat.	Mean Sat.	Med Sat.
CDS	English	5	100%	43.23%	40.00%
CDS	Spanish	29	44.83%	7.59%	6.90%
CDS	German	67	52.24%	8.31%	4.48%
Modern	English	5	100%	42.80%	40.00%
Modern	Spanish	67	43.28%	4.91%	1.49%
Modern	German	29	51.72%	5.83%	3.45%
Modern	Finnish	150	27.33%	2.46%	1.33%
Modern	Turkish	120	99.17%	4.83%	1.67%
Historic	Gothic	52	53.85%	6.31%	3.85%
Historic	Latin	113	81.42%	5.90%	2.65%

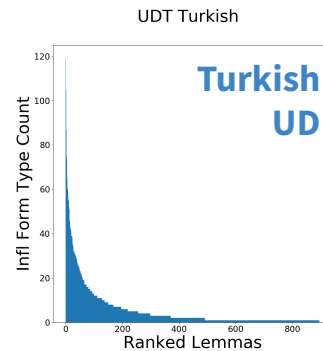
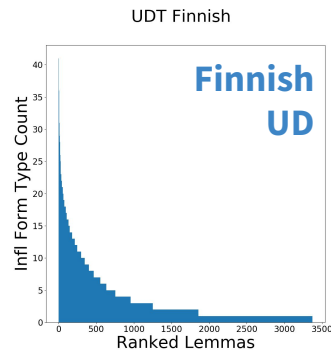
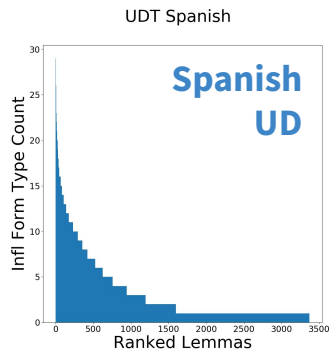
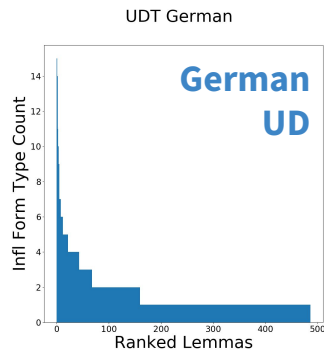
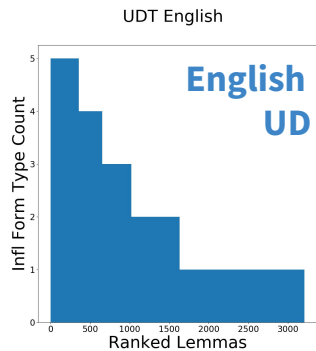
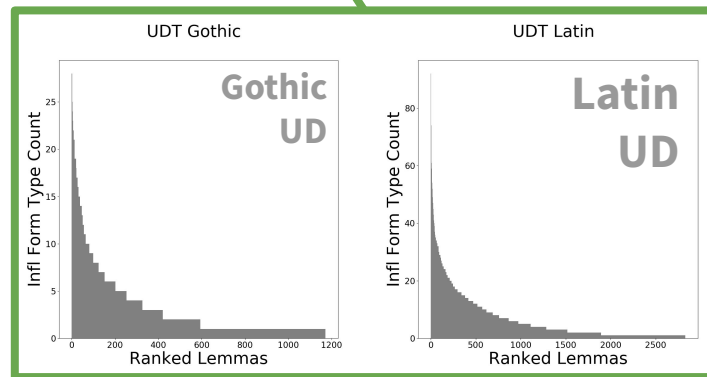
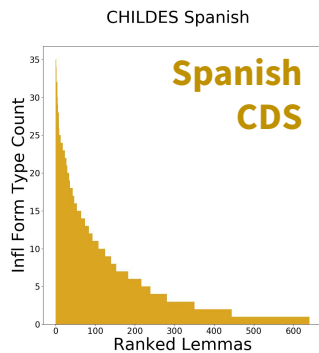
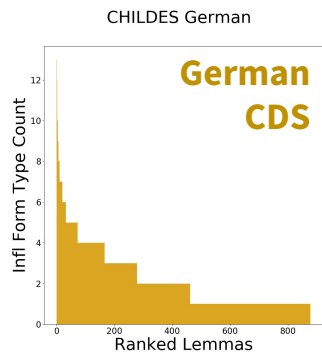
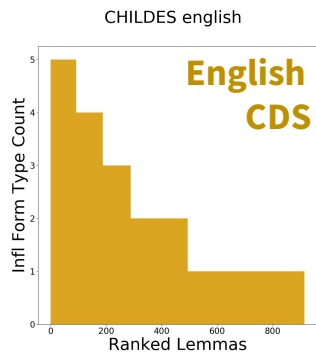
Zipfian Distributions

CDS and UD distributions correspond by language



Zipfian Distributions

Historical distributions
look like modern ones



Language Change by Language Acquisition

- Child language acquisition is one of the primary drivers of language change¹
- Not a new idea (Schleicher 1861, Paul 1880, etc)
- Children are both innovators and propagators of change
- Minor learning “errors” over successive generations → major population-level change

¹ Schleicher 1861, Paul 1880, Sweet 1899, Halle 1962, Kiparsky 1965, Andersen 1973, Baron 1977, Lightfoot 1979 *et seq*, Labov 1989, Niyogi 1996 *et seq*, Kroch 2005, Yang 2002 *et seq*, van Gelderen 2011, Cournane 2017, *inter alia*

The Paradox of Language Change

- Term coined by Niyogi & Berwick 1997
- As I see it, a central problem in the study of language change

*If children are so good at language acquisition,
why are they so bad at it?*

Transmission is not strictly linear and generational

- Children mature in communities and receive input from multiple speakers
- Young children learn sociolinguistic variables¹
- **Children attend to input from older children²** who are not linguistically mature
- **Not inconsistent with the adolescent peak³** of many continuous changes

¹ Labov 1989, Anderson 1990, ² Manly 1930, Weinreich, Labov & Herzog 1968 p 145, Roberts and Labov 1995, Labov 2001, ³ Eckert 1989, Labov 2001,

Some learning targets are unclear or absent

- One cannot acquire language from input alone due to **Poverty of the Stimulus**
- UG is proposed to render learning possible in the face of the PoS¹
- But many language specific patterns must still be acquired from the input²

Input is both richer and poorer than typically acknowledged

- Evidenced by the successes and failures of modern NLP³
- Zipfian and other long-tailed distributions for all manner of linguistic features
 - Most lexical items appear only once even in massive corpora
 - **Zipfian distributions mean sparsity is consistently worse than our intuitions about sparsity**

¹ Chomsky 1959, 1980, ² eg Baker's Paradox (Baker 1979), ³ eg the successes of distributional semantics vs the failures of coreference

Abject Poverty

Occasionally the PoS is so great that UG cannot ensure that all learners converge on the same grammar

- Forms in even moderately complex paradigms may never appear in the input¹
- Paradigmatic gaps occur when learners fail to learn a generalization for unattested input²
- Some syntactic ‘parameters’ cannot be set consistently³

¹ Chan 2008, Lignos & Yang 2017, ² Yang 2016, ³ Han et al 2007

Moving Targets

Variation is a normal and unavoidable part of acquisition

- Even in “monolingual” environments¹
- Children learn from multiple adults and each other

Change is formally inevitable²

- Given categorical representations³ and “trivial” variation
- The population composition must change over time

¹ contra Meillet, Meissel 2011, ² Niyogi & Berwick 1997, ³ Singleton & Newport 2004, Schuler et al 2017, Sneller et al 2018

What causes innovation?

“Errors” presuppose a target. Innovations need not be due to “errors”

Errors - “Blame the Child”

- The learner does not act correctly on its input “**a buggy algorithm**”
- → errors presuppose appropriate evidence and an available target

Non-errors - “Blame the Environment”

- The learner acts correctly but is dealt a bad input sample
- Even for a good algorithm, “**garbage in, garbage out**”
- Change in the face of even trivial variation

The Sibling Effect

- Why might children not overcome their overgeneralizations?

Imagine big sister Alice and little brother Bob

- Alice is currently producing innovative *ē pasts in Class IV
 - Plausible given how Class IV *ē is tenable late
 - Bob may hear these forms
- Bob is receiving both adult conservative IV pasts and Alice's
- **How does this effect Bob?**

The Sibling Effect

Can Bob identify Alice's innovation?

- Bob is likely not hear adult-produced tokens for any given low frequency Class IV verb until much later
- Since Alice is mostly consistent with adults, he cannot tell if she is innovating

Will Bob adopt Alice's innovation?

- Even young children orient toward peers
- Bob may prefer Alice's forms over his parents
- He may later learn adult forms as sociolinguistic variant doublets