

# PRODUCTIVITY OF THE SPANISH INCHOATIVE CONSTRUCTION: FROM CORPUS TO ACCEPTABILITY JUDGMENTS

Mariia Baltais (Ghent University) & Robert J. Hartsuiker (Ghent University)  
Spruik presentation at AMLaP 2022, University of York

- Syntactic productivity = a construction’s ability to attract new or existing lexical items
- Traditional corpus measures: token frequency of (co-)occurrence, type/token ratio, hapax/token ratio, etc.

How is productivity attested in corpora related to productivity “at work” in the mind of language users?

- Constructions are **extensible** beyond closed-ended corpora
  - Are corpus measures of productivity predictive of acceptability ratings?
- Grammaticality-frequency discrepancy: ratings tend to be more lenient than corpus data
  - Do ratings reflect a construction’s **extensibility**?
  - Influence of participants’ individual characteristics?

# MATERIALS

- Spanish inchoative construction [V + Prep "a" + INF]: the onset of an event
- Two sources of productivity: inchoative verb slot, infinitive slot

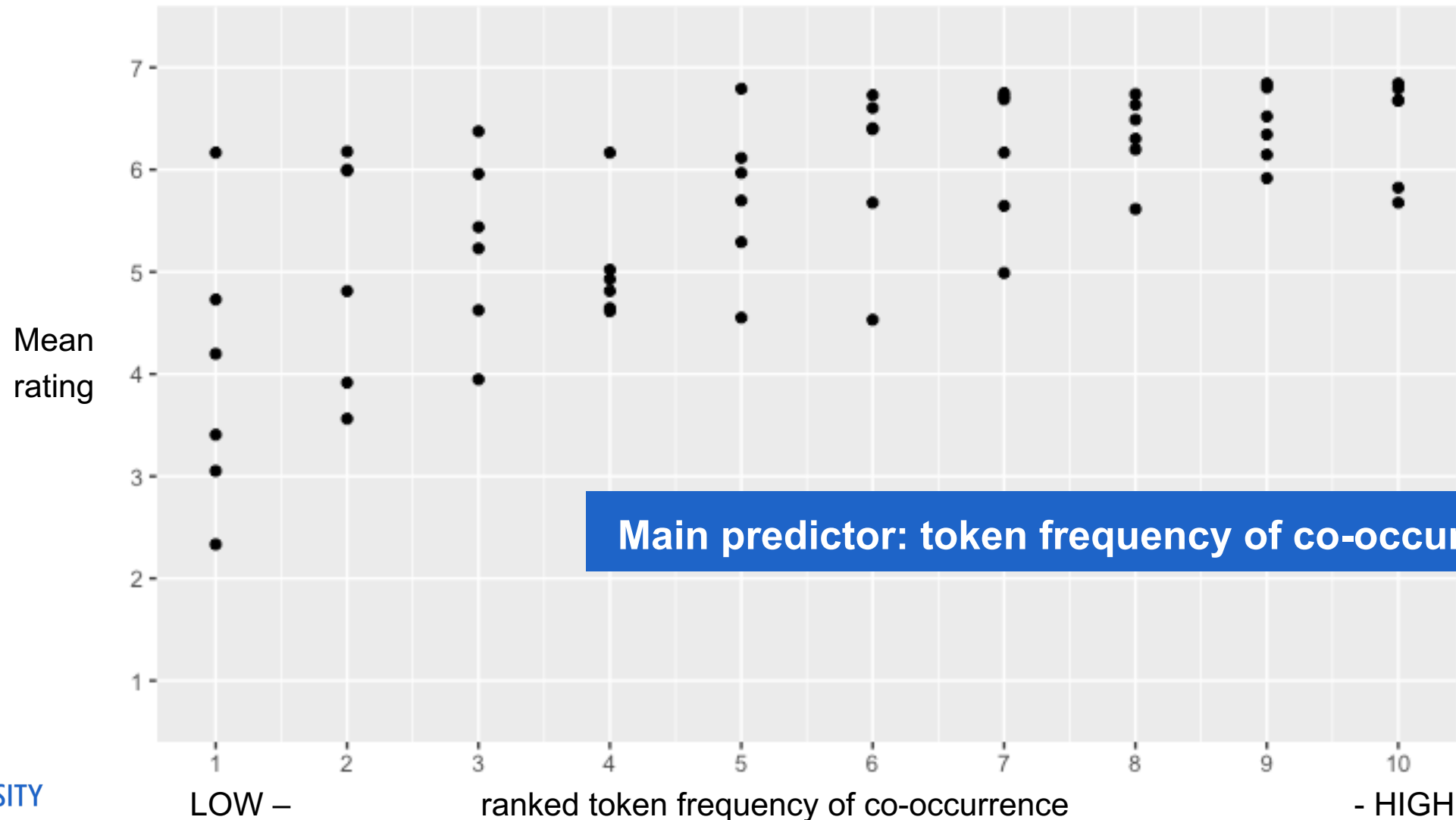
[Subj]	[V]	[Prep]	[INF]	Translation
Pedro	<b>empieza / rompe / ...</b>	a	<b>reír</b>	'Pedro begins / breaks / ... to laugh'
Pedro	<b>empieza / ?rompe</b>	a	<b>entender</b>	'Pedro begins / ?breaks to understand'

- Van Hulle & Enghels, in press: data from the Spanish Web corpus
  - 6 inchoatives with different degrees of productivity of the INF slot
  - Each combined with 10 infinitives of different token frequency of co-occurrence (incl. hapaxes and non-attested infinitives)

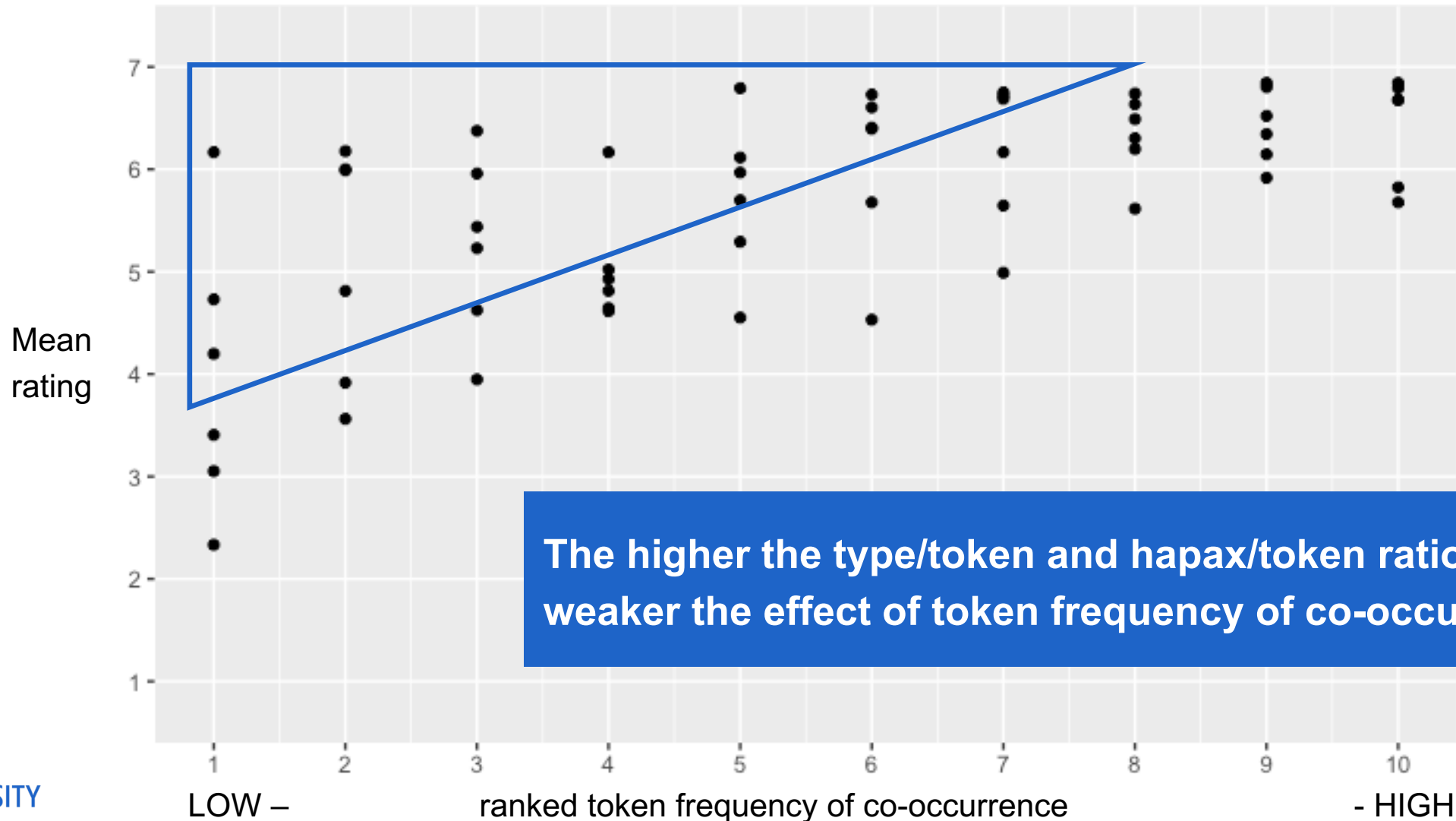
# PARTICIPANTS & PROCEDURE

- 96 native speakers of European Spanish via <https://prolific.co/>
  - 37 women, 59 men
  - Mean age: 29 y, SD: 10.4
- 60 critical sentences + 140 filler sentences = 200 in total
  - Authentic (simplified) corpus sentences
  - 21 'yes/no' comprehension questions
- 7-point Likert scale
- Sociobiographic questionnaire
- Big Five Inventory (BFI-2) for personality traits

# HIGHER FREQUENCY CORRESPONDS TO HIGHER RATINGS



# ACCEPTABILITY-FREQUENCY DISCREPANCY REFLECTS EXTENSIBILITY



# SEMANTIC COMPATIBILITY ENHANCES EXTENSIBILITY TO NON-ATTESTED ITEMS

- Subset hapaxes and non-attested INFs
- Inchoatives with high type/token and hapax/token ratios received higher ratings
- **Semantically “compatible”**\* combinations received higher ratings

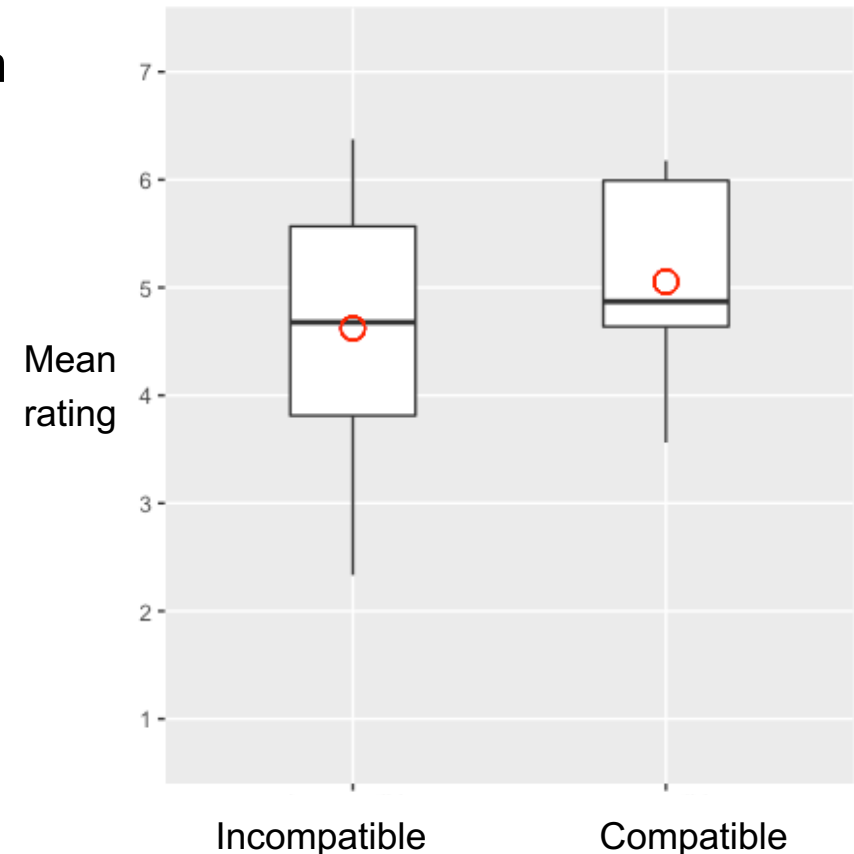
\*An infinitive is semantically compatible with the inchoative verb if it belongs to a semantic class (<http://adesse.uvigo.es/>) that is frequently attested with this inchoative in the dataset

Example compatible: class “physiology”

*echarse a respirar ‘throw oneself to breath’*

Example incompatible: class “phase”

*echarse a desarrollar ‘throw oneself to develop smth’*



# PARTICIPANT VARIABLES

Interaction with corpus measures of productivity:

- Multilingualism (more than one L2)
- Reading experience (books per year)
  - Conservative attitude towards less conventional extensions of the inchoative construction?
- No significant effects of personality traits (BFI-2)
  - Future research with more sensitive measures (e.g., eye-tracking during reading)



# DISCUSSION & CONCLUSION

*How is productivity attested in corpora related to productivity “at work” in the mind of language users?*

- Corpus measures are predictive of acceptability ratings
  - When measures reflecting different aspects of productivity are taken into account all together
- Experimental data can provide insight on **extensibility** of inchoative verbs to low-frequent and non-attested infinitives
  - Inchoatives with higher type/token and hapax/token ratios get higher ratings
  - Low-frequent infinitives are more acceptable if they are semantically “compatible” with the inchoative
- Experimental data can provide insight on the role of individual differences regarding syntactic productivity

Mariia Baltais & Robert J. Hartsuiker  
Department of Experimental Psychology

Corresponding author:  
[mariia.baltais@ugent.be](mailto:mariia.baltais@ugent.be)

“Language Productivity @ Work” project:  
<https://www.languageproductivity.ugent.be/>



Thank you!

# REFERENCES

- Baayen, R. H. (2009). 43. Corpus linguistics in morphology: morphological productivity. *Corpus linguistics. An international handbook*, 900–919.
- Barðdal, J. (2008). *Productivity: Evidence from case and argument structure in Icelandic*. Amsterdam: John Benjamins Publishing Company.
- Divjak, D. (2017). The role of lexical frequency in the acceptability of syntactic variants: Evidence from *that*-clauses in Polish. *Cognitive Science*, 41, 35–82.
- Enghels, R., & Van Hulle, S. (2018). El desarrollo de perífrasis incoativas cuasi-sinónimas: entre construccionalización y lexicalización. *ESTUDIOS DE LINGÜÍSTICA-UNIVERSIDAD DE ALICANTE-ELUA*, 32, 91–110.
- García Fernández, L. (2012). *Las perífrasis verbales en español*. Barcelona: Castalia.
- Goldberg, A. (2019). *Explain me this: Creativity, competition, and the partial productivity of constructions*. Princeton: Princeton University Press.
- Kempen, G., & Harbusch, K. (2008). Comparing linguistic judgments and corpus frequencies as windows on grammatical competence: A study of argument linearization in German clauses. In A. Steube (Ed.), *The discourse potential of underspecified structures*, 179–192. Berlin: Walter de Gruyter.
- Van Hulle, S., & Enghels, R. (in press). De Spaanse inchoatiefconstructie in beeld. Clusteranalyse als antwoord op het quasi-synonymie vraagstuk. *Handelingen – Koninklijke Zuid-Nederlandse maatschappij voor taal-en letterkunde en geschiedenis*.

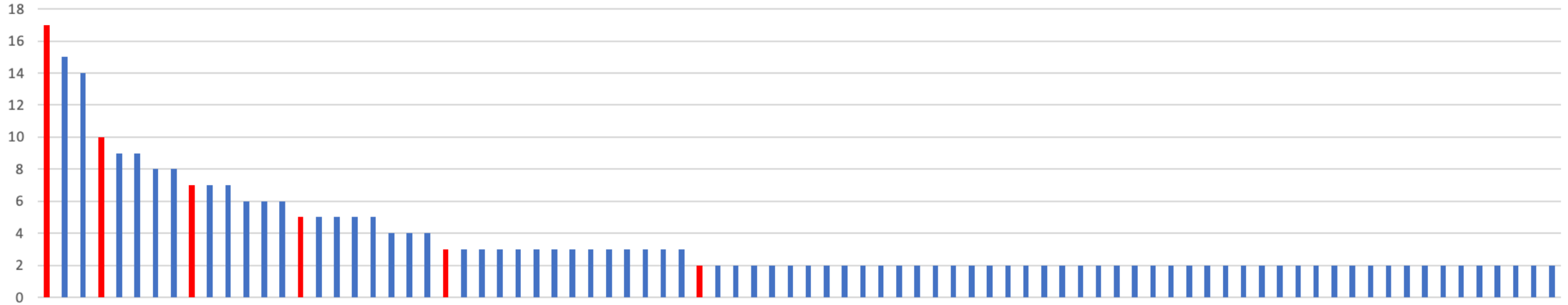
# ANNEX: CHOICE OF INCHOATIVES

- 6 inchoatives with different productivity characteristics
- Equal samples of 500 sentences each

Inchoative	Translation	Estimated token frequency in the subcorpus	Type/token ratio	Hapax/token ratio
<i>empezar</i>	'to begin, to start'	429583	0.56	0.39
<i>ponerse</i>	'to put oneself'	60728	0.36	0.24
<i>lanzarse</i>	'to throw oneself'	7476	0.43	0.27
<i>meterse</i>	'to put oneself'	1648	0.42	0.28
<i>romper</i>	'to break'	1976	0.06	0.03
<i>echarse</i>	'to throw oneself'	4242	0.03	0.01

# ANNEX: CHOICE OF INFINITIVES

*empezar* (without hapaxes)



10 infinitives per inchoative

- 6 from high to intermediate to low token frequency of co-occurrence
- 2 hapaxes: frequent and infrequent semantic class\*
- 2 non-attested: frequent and infrequent semantic class\*

\*<http://adesse.uvigo.es/>: creation, perception, displacement, physiology...

# ANNEX: ANALYSIS FULL DATA SET

Fixed effects	Estimate	SE	t-value
(Intercept)	5.64774	0.10989	51.395
Ranked token fq INCH x INF	<b>0.70366</b>	0.08626	<b>8.158*</b>
Hapax/token ratio INCH	0.13266	0.11114	1.194
Estimated token fq INCH	<b>0.28833</b>	0.11300	<b>2.551*</b>
Ranked token fq INCH x INF: Hapax/token ratio INCH	<b>-0.31554</b>	0.08483	<b>-3.720*</b>

*AR ~ Ranked token fq INCH x INF \* Hapax/token ratio + Estimated token fq INCH + (1 | item) + (1 + Ranked token fq INCH x INF \* Hapax/token ratio + Estimated token fq INCH | participant)*

- Main effect *ranked token frequency of co-occurrence*
- Main effect *estimated token fq of the inchoative in the corpus*
- Significant interaction *ranked token frequency of co-occurrence and hapax/token ratio of the inchoative*

# ANNEX: ANALYSIS HAPAXES & NON-ATTESTED

Fixed effects	Estimate	SE	<i>t-value</i>
(Intercept)	5.2037	0.2476	21.018
Hapax/token ratio INCH	<b>0.5909</b>	0.1950	<b>3.031*</b>
Semantic class (infrequent)	<b>-0.7320</b>	0.3433	<b>-2.132*</b>
Lemma fq INF	<b>-0.4823</b>	0.1683	<b>-2.866*</b>
Nu. words in sentence	-0.3783	0.1962	-1.928

*AR ~ Hapax/token ratio INCH + Semantic class + Lemma fq INF + Nu. words in sentence + (1 | item) + (1 + Hapax/token ratio INCH + Semantic class + Lemma fq INF + Nu. words in sentence | participant)*

- Main effect *hapax/token ratio of the inchoative*
- Main effect *semantic class* (incompatible vs. compatible combinations)
- Main effect *lemma frequency INF* (due to some infinitives with high lemma frequency and low ratings for the combination)

# ANNEX: ANALYSIS INDIVIDUAL VARIABLES

Fixed effects	Estimate	SE	t-value
(Intercept)	5.68461	0.14031	40.514
Token fq INCH x INF	<b>0.48396</b>	0.12726	<b>3.803*</b>
Hapax/token ratio INCH	<b>0.46705</b>	0.12758	<b>3.661*</b>
Many L2s	-0.35391	0.23977	-1.476
Books per year	0.01997	0.05581	0.358
Token fq INCH x INF: Many L2s	<b>0.23901</b>	0.08659	<b>2.760*</b>
Many L2s: Hapax/token ratio INCH	<b>0.18722</b>	0.09117	<b>2.054*</b>
Hapax/token ratio INCH: Books per year	<b>-0.05545</b>	0.02158	<b>-2.569*</b>
$AR \sim \text{Token fq INCH} \times \text{INF} * \text{Many L2s} + \text{Hapax/token ratio} * \text{Many L2s} + \text{Hapax/token ratio} * \text{Books per year} + (1   \text{item}) + (1 + \text{Token fq INCH} \times \text{INF} + \text{Hapax/token ratio}   \text{participant})$			