### Quantifying performance constraints on case-alignment typology with information theory Steven Foley %Info: https://stevenrfoley.github.io/

Conditional entropy quantifies burdens case systems place on producers and comprehenders. Entropy values for a sample of naturalistic and simulated grammars are interpolated from Georgian corpus data. Attested case patterns are simpler than expected; production constrains more than comprehension.

## 1. Information theory & grammatical complexity

• Conditional entropy (i) measures the unpredictability of one variable given a known value of another variable.

• Consider a pro-drop, scrambling SOV language with case. How hard is it to inflect NP<sub>1</sub> for case when speaking (ii)? How hard is it to parse its syntactic role while listening (iii)?

(i)  $H(Y|X) = \sum_{x \in X} \sum_{y \in Y} p(x,y) I(y|x)$  (ii)  $H(Case|Role) \equiv H(k|r) \equiv$  $I(y|x) = -\log_2 p(y|x)$ **Production Burden**  $p(y|x) = p(x,y) / p(x) \qquad (iii) \ H(Role|Case) \equiv H(r|k) \equiv$ 

• The Low Conditional Entropy Conjecture [1,2] posits a complexity ceiling on acquirable grammars.

- Inflectional paradigms can only be so irregular before the most entropic forms are regularized during language transmission.
- Research questions
- Do **syntagmatic** patterns like case alignment exhibit similarly constrained conditional entropy?
- How might H(k|r) and H(r|k) shape morphosyntactic **typology**?

## 2. Case-alignment typology & Georgian split ergativity

- Alignment refers to the way morphological categories are associated with different syntactic positions.
- Classic typology: Does case marking  $(\alpha, \beta, \gamma)$  of (iv) intransitive subjects (**S**) pattern with that of transitive subjects (A) or direct objects (**P**)?

	Α	S	Ρ
Nom-Acc	C	β	
Erg-Abs	a	3	
Tripartite	a	β	Y

- Some languages like **Georgian** [3] have very complex case alignment systems, not neatly categorizable above (iv).
- Seven syntactic roles: transitive (**A**), unergative (**Z**), unaccusative (S), and experiencer (E) subjects; patient (P) and theme (**T**) direct objects; indirect objects (**G**)
- Three-way **split**, conditioned in different environments (tenses).

v)	Macro-role		S	U		D	0	ΙΟ
	Micro-role	Α	Ζ	S	E	Р	Т	G
	Env 1		a	•	β		а	β
	Env 2	Y	Y		β	a		β
	Env 3	ļ.	3	а	β	a		(β)

**Comprehension Burden** 

## 3. Georgian case complexity calculated with corpus data

• The relative frequencies of all case–role combinations were estimated from **Georgian National Corpus** data [4].

/i)	Macro-role		S	U	D	ΙΟ		
	Micro-role	Α	Z	S	Р	Т	G	
	Env 1	37,474	15,235	38,467	14,639	43,143	26,516	26,148
	Env 2	61,928	6,572	48,717	3,596	75,825	5,356	36,980
	Env 3	9,607	574	7,947	786	12,101	1,181	3,013

vii)	Macro-role		S	U		D	0	ΙΟ
	Micro-role	Α	Ζ	S	E	Ρ	Т	G
	Case a	0.079	0.032	0.200	0	0.185	0.07	0
	Case <b>B</b>	0.020	0.001	0	0.040	0.091	0	0.129
	Case y	0.130	0.014	0	0	0	0	0

• Now calculate production & comprehension burdens of Georgian case alignment, over macro- & micro-roles. (viii)  $H(r_M|k) = 0.999$   $H(k|r_M) = 0.968$   $H(r_H|k) = 1.732$   $H(k|r_H) = 0.599$ 

# 4. Simulating a typology of case alignments

• Holding frequencies constant, corpus data were remapped to 17 attested alignments (ix) and >10k simulated ones (x).

	Α	Ζ	S	Ε	Ρ	т	G	$\mathbb{P}$
Env 1		a		β	Y	а	β	╟
Env 2		a		β	Y	а	β	╟
Env 3		a		β	Ŷ	а	β	╟
L								



Naturalistic sample: Case alignment attested in Basque, Batsbi, Cebuano, Chamorro, Chechen, Hindi, Icelandic, Inuktitut, Laz, Lezgian, Megrelian, Nez Perce, Russian, Sakha, Shipibo, Svan, Tabasaran

• Except for non-split H(k|r)s, all attested samples are less entropic than simulated ones (p < 0.05, via Welch's ind. samples *t*-tests).



Figure 1: Conditional entropy estimates for Georgian (black dots), naturalistic sample (diamonds), and simulated sample (violins/clouds). Values for split alignments are grey/black; values for non-split alignments are gold.

**References** <sup>(\*)</sup> **[1]** Ackerman, F. & R. Malouf. 2013. *Language*, 89(3). [2] Cotterell, R. et al. 2019. *Transactions of the Association for* Computational Linguistics, vol 7. [3] Harris, A. 1985. Diachronic Syntax: The Kartvelian case. [4] Gippert, J. & M. Tandashvili. 2015 In Historical Corpora, eds. J. Gippert & R. Gehrke. [5] Dell, G. & F Chang. 2014. Philosophical Transactions of the Royal Society, vol 369. [6] Futrell, R. et al. 2021. Cognitive Science, vol 44. [7] Levshina, N. 2021. Frontiers in Psychology, vol 12.

	Α	Ζ	S	Ε	Ρ	Т	G	h
1	а	β	Ŷ	а	β	Y	а	
2	Y	а	β	Y	а	β	Y	H
3	β	Y	a	β	Y	а	β	$\mathbb{H}$
								<u> </u>

Simulated sample: All 301 logically possible 3-case/7-role non-split alignments, and 10k randomly generated 2- or 3-way split alignments.

### 5. Effects of case-inventory size and number of splits

• Languages tend to have ~3 core case categories, and no more than a two-way alignment split. Is this a coincidence?



## 6. Implications for theories of case alignment

- Standard case typology is descriptive, taxonomical.
- Tentative conclusions
- sample are **verb-initial**!



• Question for 'theoretical typology': how do grammatical factors contribute to complexity, independent of alignment?

• To explore this question: 16 new simulated typologies of 5k alignments, varying in numbers of cases and splits

Information theory provides an **explanatory foothold**.

• Sentence-processing theories like the P-Chain [5] might explain a complexity ceiling on case: more entropic grammars are harder to use, making acquisition channels noisier [cf. 6].

• Attested case-alignment patterns are **less entropic** than they could be (unsurprising — how would very complex ones arise?) • Case is more streamlined for production than comprehension, especially when calculated over syntactic micro-roles. Systems with more cases benefit comprehension but impede production; **more splits** make everything harder. • But **2-case** languages are the hardest to parse — and rarer? • Limitations of this approach: highly **abstract**, doesn't

account for other cues like word order or animacy [cf. 7]. • Tantalizing observation: the most entropic languages in the