

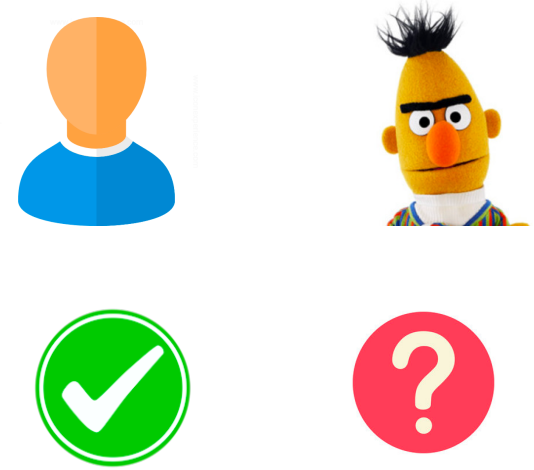


Les Fondation Modèles comprennent **presque** les quantificateurs.

## Research Question

Can foundation models understand (generalized) quantifiers like humans?

Some birds can fly.  $\rightarrow$   $X\text{-}Y\%$  ( $0 < X < Y < 100$ ) birds can fly.



## NLI Limitations

- Implicit percentage values of quantifiers.
- Sentence-level relation nature; impacts of linguistic and social clues.
- Deficiencies in ambiguous premises and quantitative reasoning.

## Task

- Splitting  $[0, 1]$  in  $W$  with  $\beta$ :  $W_{\beta=0.05} = \{0.5\%, 10\% \dots\}$
- A model receives a quantified sentence and outputs  $[p_1, p_2] \in W_{\beta}$  where the predicate in the quantified sentence holds true.

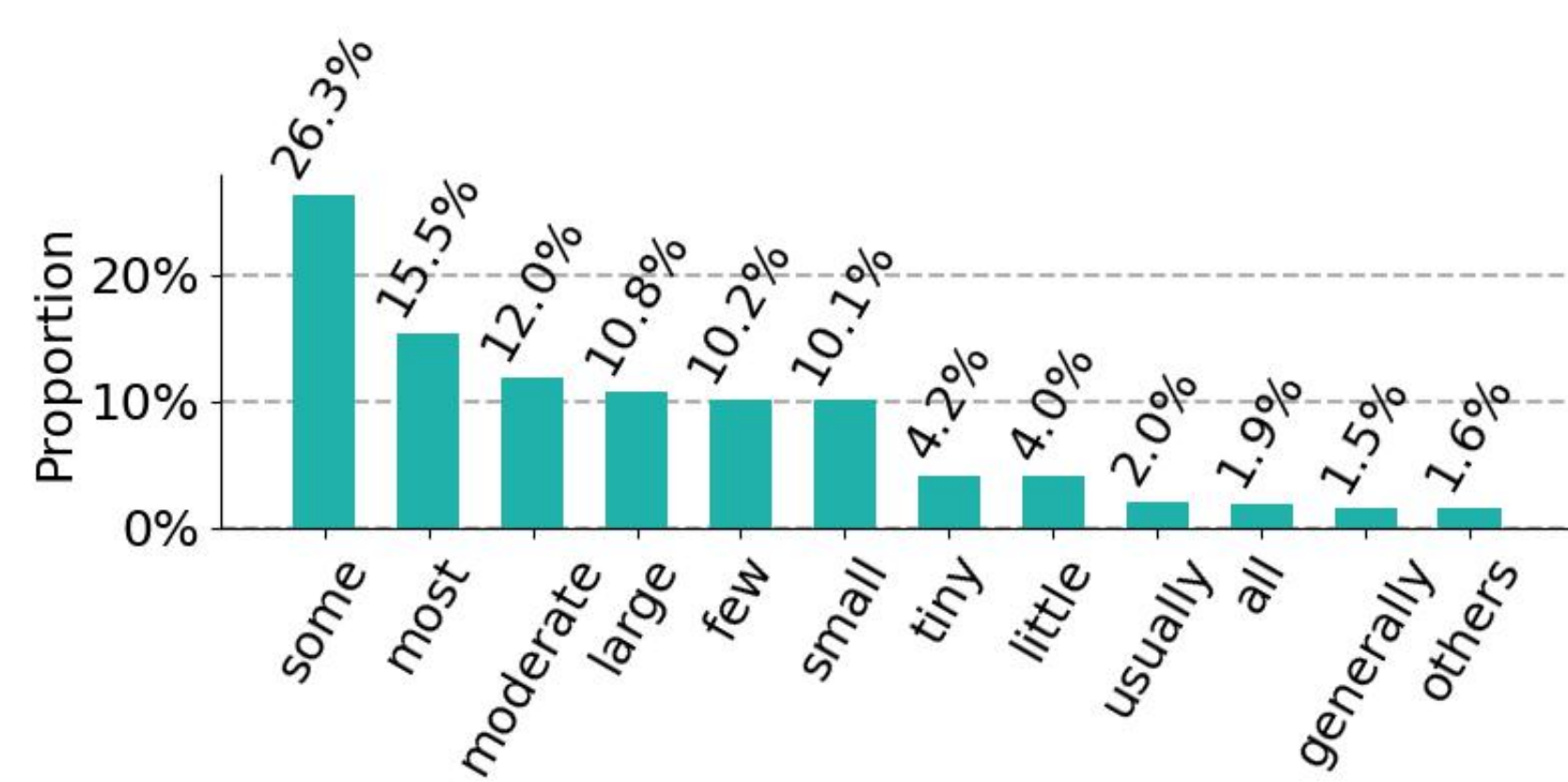
## PRESQUE: Pragmatic Reasoning for Semantics of Quantifiers



- NLI backbone
- RSA adoption
- No training, no bias

## New Dataset: QuRe

- Existing datasets (e.g. HVD):
  - Limited quantifiers
  - **No golden percentage scope**
  - e.g. All rocks have minerals.
- Ours: QuRe
  - **More** generalized quantifiers

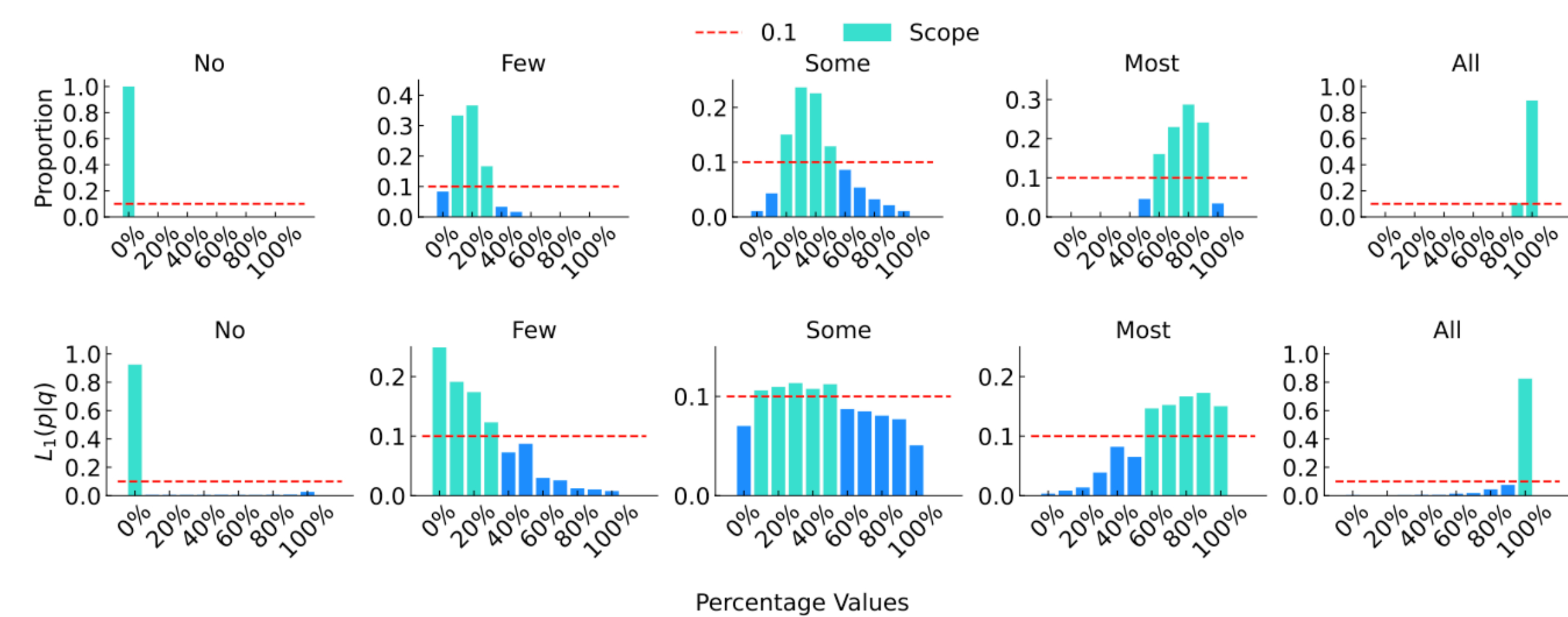


- **Specificity** levels
  - Difficulty to reason percentage scope from the non-quantification context.
  - Full/Partially/Indeterminable
- **Golden percentage scope**
- Sentence **topics**

[WIKI ENTITY] ORIGINAL SENTENCE	[SPECIFICITY, EXPRESSION] QuRe SENTENCE	TOPICS
[Human] Most humans (61%) live in Asia; the remainder live in the Americas (14%), Africa (14%), Europe (11%), and Oceania (0.5%). Within the last century, humans have explored challenging environments such as Antarctica, the deep sea, and outer space.	[Fully, 0.11] Most humans (61%) live in Asia; the remainder live in the Americas (14%), Africa (14%), some Europe, and Oceania (0.5%). Within the last century, humans have explored challenging environments such as Antarctica, the deep sea, and outer space.	population continents exploration
[List of blade materials] Prior to 2002, INFI contained 0.5% Carbon, 0.74% Nitrogen, about 1% Cobalt, and about 0.1% Nickel.	[Partially, 0.005] Prior to 2002, INFI contained tiny levels of Carbon, 0.74% Nitrogen, about 1% Cobalt, and about 0.1% Nickel.	chemical composition INFI elements
[List of blade materials] In order for a steel to be considered stainless it must have a Chromium content of at least 10.5%.	[Indeterminable, >= 0.105] In order for a steel to be considered stainless it must have some Chromium content.	steel metallurgy composition

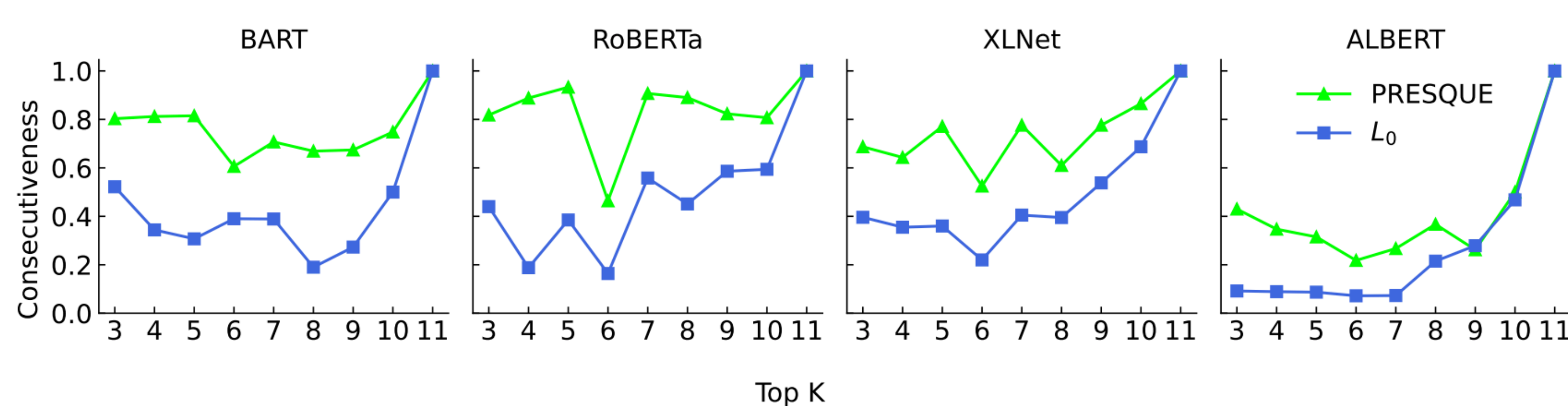
## Human & PRESQUE Perception

Human defined percentage scopes of quantifiers are similar to PRESQUE.



## Consecutiveness

- If top 3 percentage predictions are
  - $\{10\%, 20\%, 30\%\}$ , gives a consecutive scope 10%-30%. **Consistent**
  - $\{10\%, 30\%, 50\%\}$ , no consecutive scope. **Not consistent**
- PRESQUE predictions has higher consecutiveness ratio than  $L_0$ .



## Beyond Direct NLI Interpretation

- Pragmatic theory (Grice 1975)
- Rational Speech Act (RSA, Goodman and Frank 2016)
  - Word state  $W$  and utterances  $U$ .
    - $W$ : percentage value set  $\{0, 10\%, 20\% \dots\}$
    - $U$ : quantifier set  $\{\text{no, few, some} \dots\}$
  - The mental states of Listener  $L$  and speaker  $S$  are iteratively modeled with a Bayesian approach.

$$\begin{array}{lcl} L_1 & \text{pragmatic listener} & P_{L_1}(s|u) \propto P_{S_1}(u|s) \cdot P(s) \\ \downarrow & & \\ S_1 & \text{pragmatic speaker} & P_{S_1}(u|s) \propto \exp(\alpha U_{S_1}(u;s)) \\ \downarrow & & \\ L_0 & \text{literal listener} & P_{L_0}(s|u) \propto \llbracket u \rrbracket(s) \cdot P(s) \end{array}$$

## PRESQUE ( $L_1$ )

- Premise  $\tilde{p}$ : All airplanes have engines.
- Hypothesis  $\tilde{h}$ : 90% airplanes have engines.
- $p$ : percentage values.
- $q$ : quantifiers
- Entailment( $\cdot$ ) comes from the NLI model.

- Literal listener

$$L_0(p|q) \propto \text{Entailment}(\tilde{p}, \tilde{h})$$

- Pragmatic speaker

$$S_0(q|p) \propto \text{Entailment}(\tilde{h}, \tilde{p})$$

- Pragmatic listener

$$L_1(p|q) \propto S_0(q|p)P(p) \quad P(p) = \sum_{q \in \mathcal{U}} P(p|q)P(q)$$

## Result (QuRe)

- PRESQUE generally performs better than  $L_0$  among all

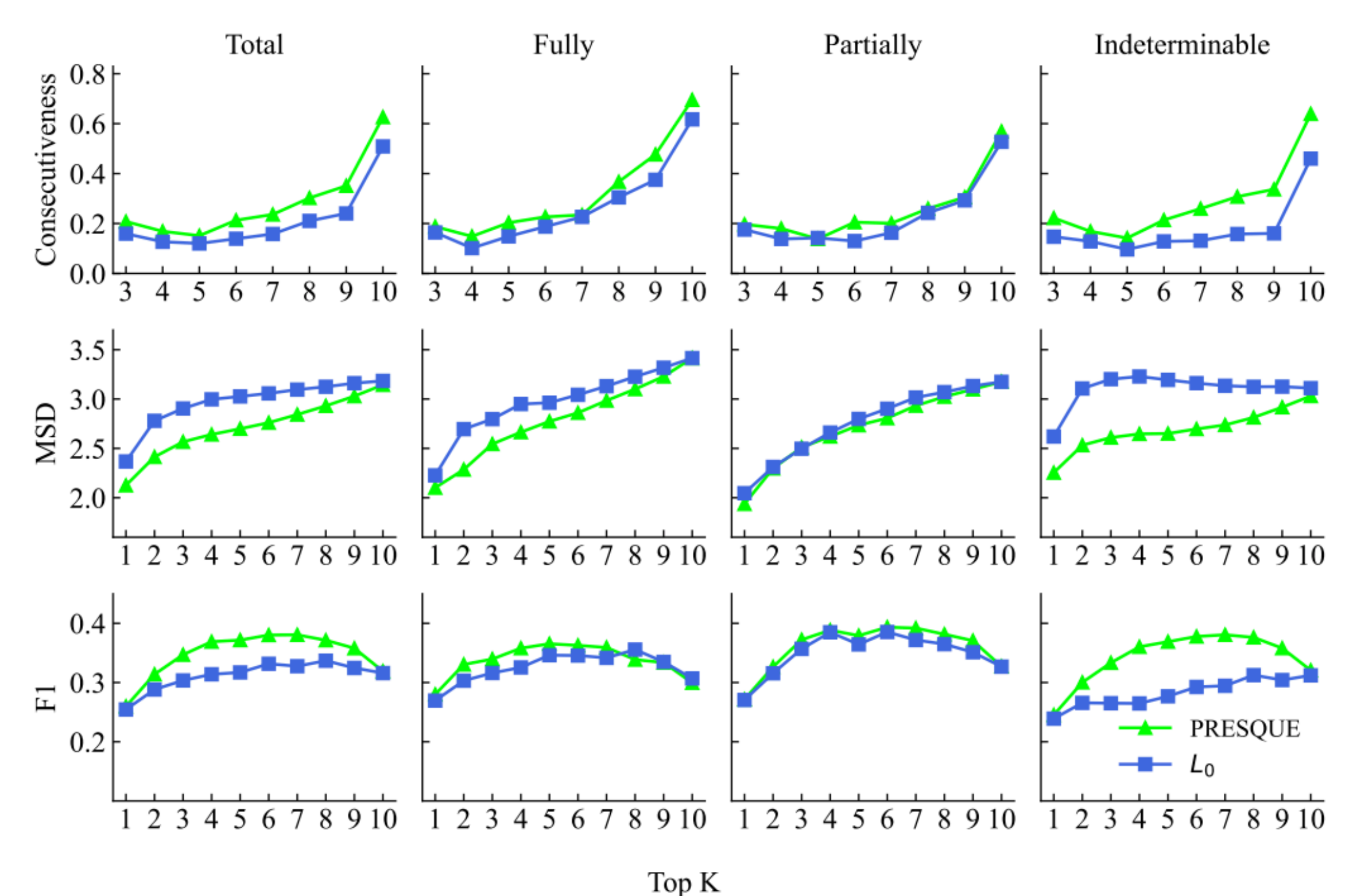
SPECIFICITY	HIT@1 $\uparrow$			MRR $\uparrow$			CROSSENTROPY $\downarrow$			F1@{1, 5} $\uparrow$		
	Rnd.	$L_0$	$L_1$	Rnd.	$L_0$	$L_1$	Rnd.	$L_0$	$L_1$	Rnd.	$L_0$	$L_1$
Fully	4.1	27.3	<b>29.7</b>	12.3	22.1	<b>24.3</b>	6.44	<b>5.64</b>	5.74	2.8/8.6	19.5/24.3	<b>21.5/26.5</b>
Partial	8.2	26.4	<b>28.5</b>	11.6	21.2	<b>21.7</b>	7.78	<b>6.99</b>	7.06	4.3/8.3	16.9/25.9	<b>18.3/27.3</b>
Indeterminable	9.7	21.4	21.4	12.5	18.1	<b>22.7</b>	7.76	7.20	<b>6.69</b>	5.3/10.1	<b>14.9/18.2</b>	14.8/25.6
Total	7.9	24.0	<b>25.1</b>	11.8	19.8	<b>22.7</b>	7.47	6.86	<b>6.78</b>	4.4/9.3	16.3/21.7	<b>17.1/26.3</b>

## Examples

[GS.] SENTENCE <sub>Q</sub> / [SPC.] SENTENCE <sub>P</sub>	PRIMARY SCOPE	MRR	F1@5	CE
[F] In 57 separate fights, one loss was observed to Neope goschkevitschii, giving V. mandarinia a large winning rate.	$L_0$ : 5%-20%	0.11	0.00	7.67
[95%-100%] In 57 separate fights, one loss was observed to Neope goschkevitschii, giving V. mandarinia a win rate of 98.3%.	$L_1$ : 85%-100%	<b>0.67</b>	<b>0.67</b>	<b>3.52</b>
[P] From 4 locations in different parts of Europe, a large number had clutch size of 2, 41% had size of 3, clutches of 1 and 4 each constituted about 8%.	$L_0$ : 30%-40%	0.22	0.40	6.29
[40%-45%] From 4 locations in different parts of Europe, 43% had clutch size of 2, 41% had size of 3, clutches of 1 and 4 each constituted about 8%.	$L_1$ : 30%-45%	<b>0.33</b>	<b>0.67</b>	<b>4.92</b>
[I] It is typically made from rye bread, usually known as black bread, and is not classified as an alcoholic beverage in Poland, as its alcohol content usually is very little.	$L_0$ : 60%-70%	0.06	0.00	6.97
[0-5%] It is typically made from rye bread, usually known as black bread, and is not classified as an alcoholic beverage in Poland, as its alcohol content usually ranges from 0% to 2%.	$L_1$ : 0%-5%	<b>0.33</b>	<b>1.00</b>	<b>4.16</b>

## Consecutiveness & Distance Metrics

- MSD: the minimal distance between the consecutive scope of top K percentage predictions and the golden percentage scope.
- F1: the span overlap between the consecutive scope of top K percentage predictions and the golden percentage scope.
- PRESQUE has higher consecutiveness and lower MSD



## Pragmatic Reasoning Unlocks Quantifier Semantics for Foundation Models

Yiyuan Li, Rakesh R Menon, Sayan Ghosh, Shashank Srivastava  
[yiyanli@cs.unc.edu](mailto:yiyanli@cs.unc.edu)



Paper

