



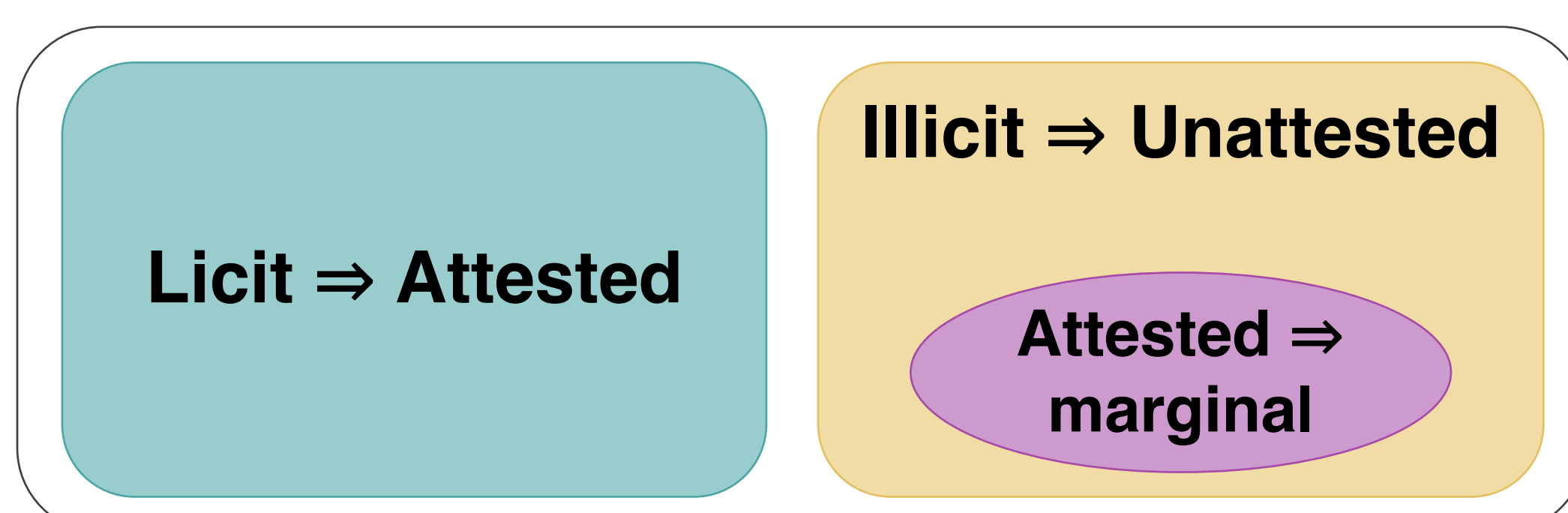
## MARGINAL SEQUENCES IN PHONOTACTIC THEORY

### How are ATTESTATION and LICITNESS related?

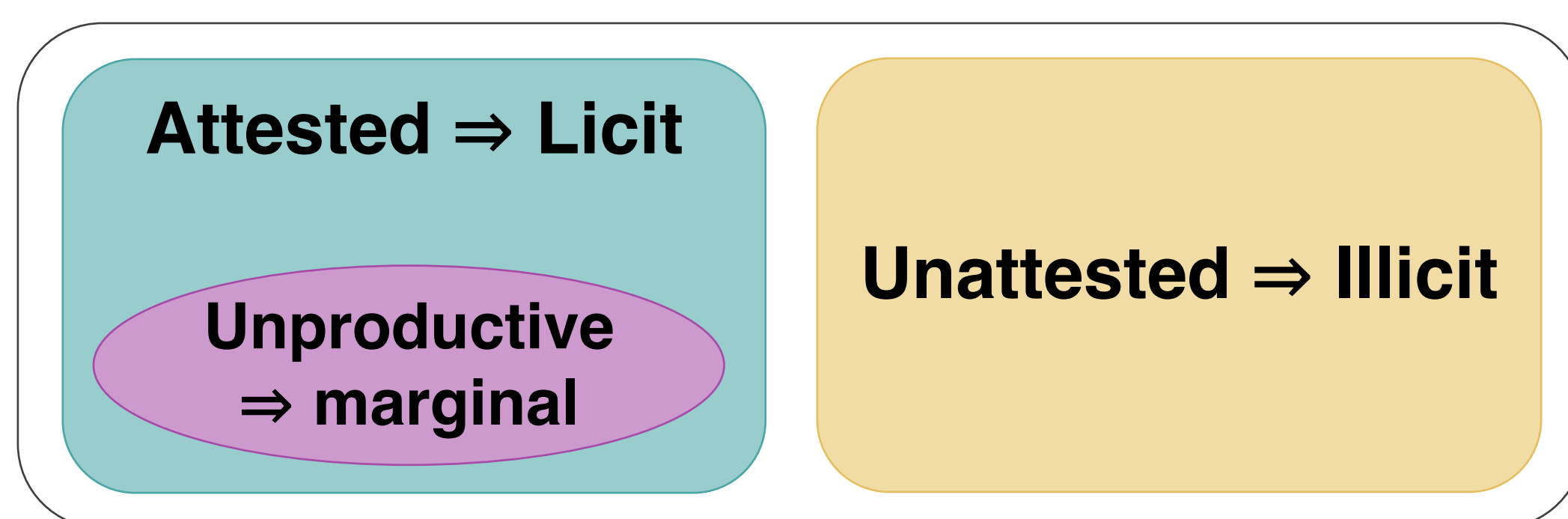
- **ATTESTED** subsequences are generally **LICIT**
- **UNATTESTED** subsequences are generally **ILLICIT**

### Where do MARGINAL FORMS fit in?

- Previous approaches: **ILLICIT BUT ATTESTED** (Hyman 1975)



- Our approach: **ATTESTED BUT UNPRODUCTIVE**



## EVIDENCE FOR OUR MODEL

- **BORROWINGS:** not repaired

	Spanish	Japanese	English
German: /pfɪtse/	/fajser/	/φaidza/	/faɪzɪ/
Italian: /spagetti/	/espageti/	/swpagetti/	/spəgeti/
Greek: /sfɪŋks/	/esfinxe/	/swφinkɯsw/	/sfɪŋks/
Greek: /sfaira/	/esfera/	(swφia)	/sfɪə/

- **NEW WORDS:** may contain marginal sequences



spheal



dwebble

- **PRODUCTION & PERCEPTION ERRORS**

- Speakers **struggle to produce illicit sequences**
- **97% production accuracy** on /#sC/ sequences by English speakers
  - $C \in \{f, p, t, k, m, n\}$  (Davidson 2006)

## FORMALIZING MARGINAL VS. LICIT WITH THE TSP

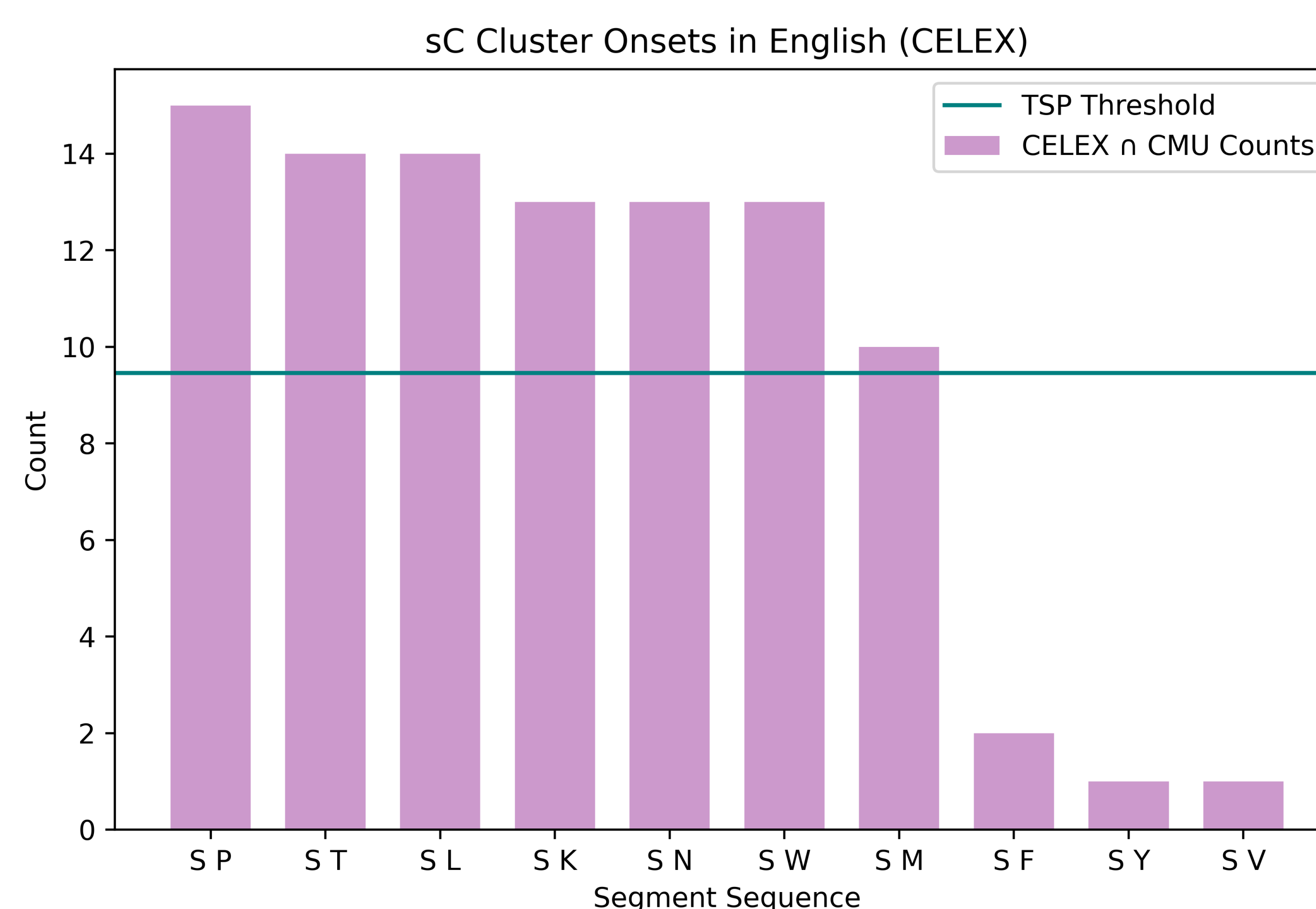
- **LICIT VS. MARGINAL AS A DIFFERENCE IN PRODUCTIVITY**

- **LICIT ONSETS/CODAS:** occur with a **sufficiently diverse set of nuclei**
  - Occurs with **most nuclei**  $\Rightarrow$  can occur with all
- **MARGINAL ONSETS/CODAS:** can occur with only a few, **lexicalized** nuclei
  - Occurs with a **few nuclei**  $\Rightarrow$  memorize those nuclei

- **THE TOLERANCE PRINCIPLE (TSP, YANG 2016):**

- In a language with  $N$  possible nuclei, an attested onset/coda is **LICIT** if it occurs with at least  $M$  of these possible nuclei and:

$$N - M \leq \theta_N = \frac{N}{\ln N}$$



## DATA

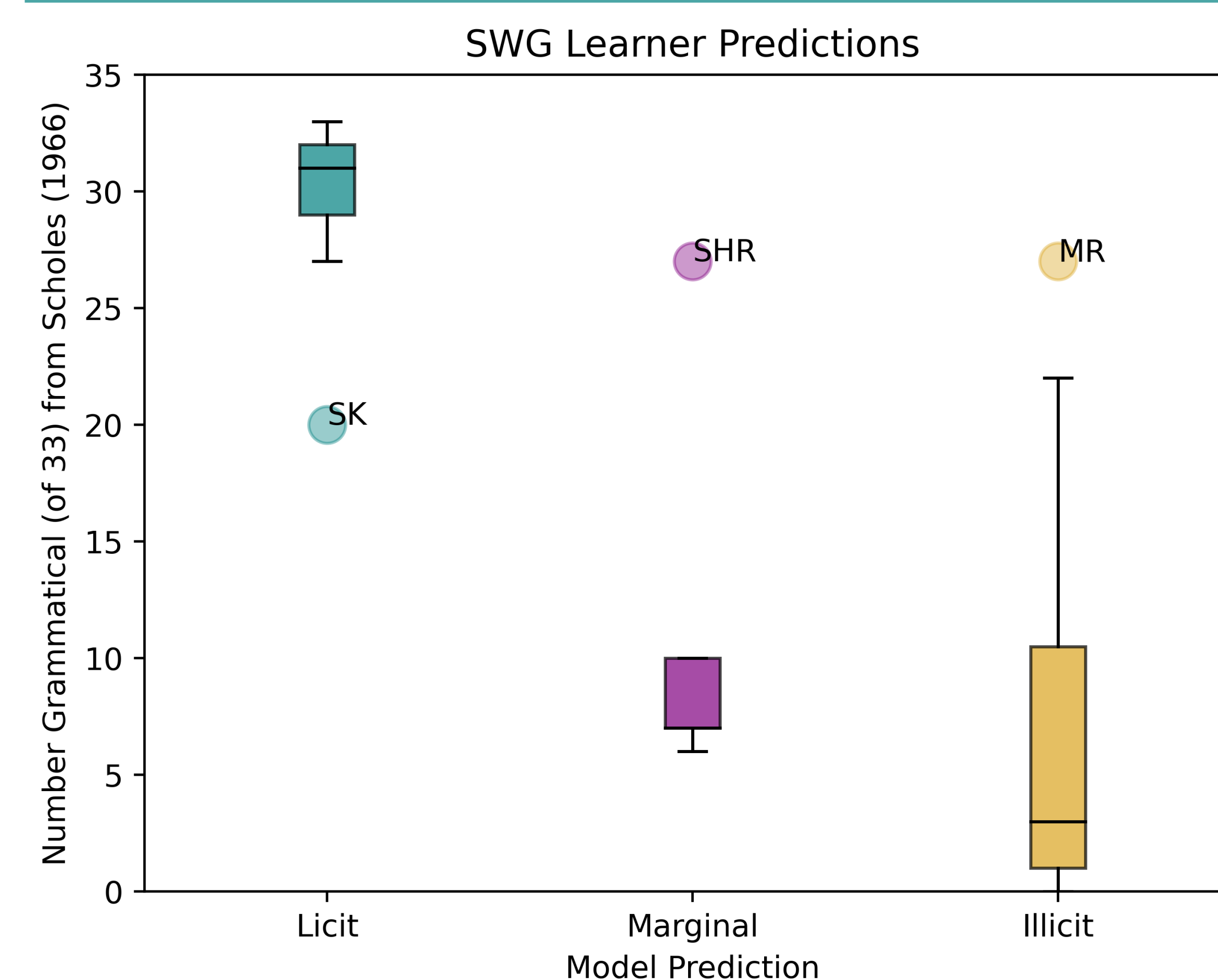
- **TRAINING:**

- **CELEX n CMU PRONOUNCING DICTIONARY:** ~41k words
- Syllabify and extract syllable constituents (Gorman 2013)
- Phonological Features from Hayes & Wilson 2008

- **JUDGMENTS:**

- **SCHOLES:** complex onsets in monosyllabic nonce words
- Binary decisions by 33 seventh graders

## RESULTS



	Attestation	SWG	H&W
Pearson's $r$	0.78	0.86	0.84
Spearman's TR $\rho$	0.74	0.78	0.79
Goodman-Kruskal $\gamma$	0.89	0.89	0.65
Kendall's $\tau_b$	0.62	0.66	0.61

## MODEL: SEQUENCE-WISE GENERALIZATION LEARNER (SWG)

- **MOTIVATION & ASSUMPTIONS:**

- Phonotactic knowledge **represented over syllables**
- Representations initially **featurally-underspecified** during acquisition

We present a **SYLLABLE-BASED** computational model that learns a **POSITIVE PHONOTACTIC GRAMMAR** categorizing forms as **LICIT, MARGINAL, OR ILLICIT**.

- **LEARNING ALGORITHM:** recursive, feature-based subdivision to learn phonotactics as **increasingly-specific sequences of feature sets**

- At each step, intersect all sequences in current input to give underspecified sequence  $S$
- If **sufficiently many sequences** matching  $S$  are licit, **add  $S$  to set of licit sequences**
- Otherwise, **subdivide the input** based on the most frequent feature at the index in the string with the greatest difference between  $N$  and  $M$ , and recurse
- If **no generalization & no more features to subdivide on**, then  $S$  is **marginal**

### ACKNOWLEDGEMENTS:

I am grateful to Jeff Heinz, Jordan Kodner, Charles Yang, Scott Nelson, Salam Khalifa, Felix Fonseca, Kyle Gorman, and Huteng Dai for helpful discussion. This work was supported by the Institute for Advanced Computational Science (IACS) Graduate Research Fellowship and the National Science Foundation (NSF) Graduate Research Fellowship Program under NSF Grant No. 2234683. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author and do not necessarily reflect the views of the IACS or the NSF.

## FUTURE WORK

- **DEVELOPMENTAL IMPLICATIONS**

- Model predicts **initial stage of conservatism**
- Must **accumulate sufficient evidence** for licitness

- **FURTHER COMPARISONS**

- Human judgments on English & **other languages**
- Comparison to H&W and **other models**

- How can we learn **SYLLABLE CONTACT CONSTRAINTS** in this framework?

- How does SWG fare on languages with **SMALLER VOWEL SPACES?**

- Prediction: **more onsets/codas will pass TSP** and be licit because  $N$  will be smaller

### SELECTED REFERENCES:

- Hayes & Wilson 2008. A maximum entropy model of phonotactics and phonotactic learning. *Linguistic Inquiry*.
- Davidson 2006. Phonology, phonetics, or frequency: Influences on the production of non-native sequences. *Journal of Phonetics*.
- Gorman 2013. *Generative Phonotactics*. UPenn Dissertation.
- Hyman 1975. *Phonology: Theory and Analysis*. Harcourt Press.
- Kabak & Idsardi 2007. Perceptual distortions in the adaptation of English consonant clusters: Syllable structure or consonantal contact constraints? *Language and Speech*.
- Yang 2016. *The Price of Linguistic Productivity: How Children Learn to Break the Rules of Language*. MIT Press.