

# Byung-Doh Oh

✉ oh.b@nyu.edu  
🌐 Professional Website  
🐙 GitHub Profile  
🎓 Google Scholar

## RESEARCH INTERESTS

---

- Computational Linguistics; Psycholinguistics; Cognitive Modeling; Natural Language Processing; Large Language Models; Neural Network Interpretability
- My work aims to advance our understanding of real-time language processing in humans and machines by drawing on methods from computational psycholinguistics and machine learning.

## EDUCATION

---

<b>The Ohio State University</b>	Columbus, OH
Ph.D. in Linguistics	2024
<ul style="list-style-type: none"><li>• Dissertation: <i>Empirical Shortcomings of Transformer-Based Large Language Models as Expectation-Based Models of Human Sentence Processing</i></li><li>• Advisor: William Schuler</li><li>• Committee: Michael White, Micha Elsner, Tal Linzen</li></ul>	
<b>Seoul National University</b>	Seoul, Korea
M.A. in English Language Education	2018
<ul style="list-style-type: none"><li>• Thesis: <i>Exploring English Online Research and Comprehension Strategies of Korean College Students</i></li><li>• Advisor: Youngsoon So</li></ul>	
<b>Seoul National University</b>	Seoul, Korea
B.A. in English Language Education, <i>summa cum laude</i>	2016
<ul style="list-style-type: none"><li>• Thesis: <i>A Study on the Assessment Techniques of Language MOOCs</i></li><li>• Advisor: Youngsoon So</li></ul>	

## EMPLOYMENT

---

<b>Assistant Professor/Faculty Fellow</b>	Sept. 2024 –
Center for Data Science, New York University	New York, NY
<b>Natural Language &amp; Speech Processing Research Intern</b>	May 2021 – Aug. 2021
Tencent AI Lab	Bellevue, WA (remote)
<ul style="list-style-type: none"><li>• Mentor: Lifeng Jin</li></ul>	

## AWARDS AND HONORS

---

<b>CDS Faculty Fellowship</b>	2024 – 2026
Center for Data Science, New York University	
<b>Selection as DARPA Riser</b>	2022
Defense Advanced Research Projects Agency	
<b>Fulbright Graduate Study Award</b>	2018 – 2020
Bureau of Educational and Cultural Affairs	
<b>Chunjae Education Group Scholarship</b>	2015
Chunjae Education Group	
<b>College of Education Alumni Association Scholarship</b>	2011
Seoul National University	

- Under review **Byung-Doh Oh** and William Schuler. Dissociable frequency effects attenuate as large language model surprisal predictors improve. *Journal of Memory and Language*.
- 2024 **Byung-Doh Oh** and William Schuler. The impact of token granularity on the predictive power of language model surprisal. *arXiv preprint*, arXiv:2412.11940.

---

PEER-REVIEWED ARTICLES AND PROCEEDINGS PAPERS

---

- 2025 Christian Clark, **Byung-Doh Oh**, and William Schuler. Linear recency bias during training improves Transformers' fit to reading times. In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 7735–7747.
- 2024 **Byung-Doh Oh** and William Schuler. Leading whitespaces of language models' subword vocabulary pose a confound for calculating word probabilities. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 3464–3472.
- 2024 **Byung-Doh Oh**, Shisen Yue, and William Schuler. Frequency explains the inverse correlation of large language models' size, training data amount, and surprisal's fit to reading times. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics*, pages 2644–2663.
- 2023 **Byung-Doh Oh** and William Schuler. Transformer-based language model surprisal predicts human reading times best with about two billion training tokens. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 1915–1921.
- 2023 **Byung-Doh Oh** and William Schuler. Token-wise decomposition of autoregressive language model hidden states for analyzing model predictions. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics*, pages 10105–10117.
- 2023 **Byung-Doh Oh** and William Schuler. Why does surprisal from larger Transformer-based language models provide a poorer fit to human reading times? *Transactions of the Association for Computational Linguistics*, 11:336–350.
- 2022 **Byung-Doh Oh** and William Schuler. Entropy- and distance-based predictors from GPT-2 attention patterns predict reading times over and above GPT-2 surprisal. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 9324–9334.
- 2022 **Byung-Doh Oh**, Christian Clark, and William Schuler. Comparison of structural parsers and neural language models as surprisal estimators. *Frontiers in Artificial Intelligence*, 5:777963.
- 2021 Lifeng Jin, **Byung-Doh Oh**, and William Schuler. Character-based PCFG induction for modeling the syntactic acquisition of morphologically rich languages. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 4367–4378.
- 2021 Evan Jaffe, **Byung-Doh Oh**, and William Schuler. Coreference-aware surprisal predicts brain response. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 3351–3356.
- 2021 **Byung-Doh Oh**, Christian Clark, and William Schuler. Surprisal estimators for human reading times need character models. In *Proceedings of the Joint Conference of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing*, pages 3746–3757.
- 2021 **Byung-Doh Oh** and William Schuler. Contributions of propositional content and syntactic category information in sentence processing. In *Proceedings of the 11th Workshop on Cognitive Modeling and Computational Linguistics*, pages 241–250.

- 2021 **Byung-Doh Oh**. Team Ohio State at CMCL 2021 shared task: Fine-tuned RoBERTa for eye-tracking data prediction. In *Proceedings of the 11th Workshop on Cognitive Modeling and Computational Linguistics*, pages 97–101.
- 2019 **Byung-Doh Oh\***, Pranav Maneriker\*, and Nanjiang Jiang\*. THOMAS: The hegemonic OSU morphological analyzer using seq2seq. In *Proceedings of the 16th Workshop on Computational Research in Phonetics, Phonology, and Morphology*, pages 80–86.
- 2019 Micha Elsner, Andrea D. Sims, Alexander Erdmann, Antonio Hernandez, Evan Jaffe, Lifeng Jin, Martha Booker Johnson, Shuan Karim, David L. King, Luana Lamberti Nunes, **Byung-Doh Oh**, Nathan Rasmussen, Cory Shain, Stephanie Antetomaso, Kendra V. Dickinson, Noah Diewald, Michelle McKenzie, and Symon Stevens-Guille. Modeling morphological learning, typology, and change: What can the neural sequence-to-sequence framework contribute? *Journal of Language Modelling*, 7(1):53–98.
- 2018 **Byung-Doh Oh** and Youngsoon So. Exploring English online research and comprehension strategies of Korean college students. *English Teaching*, 73(3):53–76.
- 2017 **Byung-Doh Oh**. Predicting L2 writing proficiency with computational indices based on n-grams. *Foreign Language Education Research*, 21:1–20.

#### INVITED TALKS

---

- 2025 **Byung-Doh Oh**. *Unveiling the language processing of humans and machines*. Invited talk at the Linguistics and Multilingual Studies Program, Nanyang Technological University, Singapore.
- 2024 **Byung-Doh Oh**. *What can linguistic data tell us about the predictions of large language models?* Invited talk at the Department of Computer Science and Engineering and Graduate School of Artificial Intelligence, POSTECH, Pohang, Korea.
- 2024 **Byung-Doh Oh**. *The bigger-is-worse effects of model size and training data of large language model surprisal on human reading times*. Invited talk in the Distinguished Speakers in Language Science Colloquium Series, Saarland University, Saarbrücken, Germany.
- 2023 **Byung-Doh Oh**. *The bigger-is-worse effects of model size and training data of large language model surprisal on human reading times*. Invited talk at the Center for Data Science, New York University, New York, NY.
- 2022 **Byung-Doh Oh**. *Computational models of sentence processing and syntactic acquisition*. Invited talk at the Department of English, Dongguk University, Seoul, Korea (online).

#### CONFERENCE PRESENTATIONS

---

- 2025 **Byung-Doh Oh** and William Schuler. *Word frequency modulates the effects of model size and training data amount on language model surprisal*. Poster presented at the 38th Annual Conference on Human Sentence Processing, College Park, MD.
- 2025 **Byung-Doh Oh** and William Schuler. *Correcting language model word probabilities reveals a greater divergence between surprisal and human reading times*. Poster presented at the 38th Annual Conference on Human Sentence Processing, College Park, MD.
- 2023 **Byung-Doh Oh** and William Schuler. *On the bigger-is-worse nature of pre-trained language model surprisal*. Poster presented at the 36th Annual Conference on Human Sentence Processing, Pittsburgh, PA.
- 2023 **Byung-Doh Oh** and William Schuler. *Memory-based predictors from GPT-2 attention predict reading times over surprisal*. Poster presented at the 36th Annual Conference on Human Sentence Processing, Pittsburgh, PA.

2022	<b>Byung-Doh Oh.</b> <i>Unified unsupervised grammar induction for typologically diverse languages.</i> Poster presented at the DARPA Risers program, Columbus, OH.
2021	<b>Byung-Doh Oh</b> , Christian Clark, and William Schuler. <i>Comparison of structural and neural language models as surprisal estimators.</i> Short talk presented at the 34th Annual CUNY Conference on Human Sentence Processing, Philadelphia, PA (online).
2021	<b>Byung-Doh Oh</b> and William Schuler. <i>Contributions of propositional content and syntactic categories in sentence processing.</i> Short talk presented at the 34th Annual CUNY Conference on Human Sentence Processing, Philadelphia, PA (online).
2019	Evan Jaffe and <b>Byung-Doh Oh.</b> <i>The role of learnability in morphological change: A computational approach.</i> Talk presented at the Fourth American International Morphology Meeting, Stony Brook, NY.

## TEACHING

---

2025	<b>Invited Lecture</b> , HG2051: Language and the Computer, Nanyang Technological University (Instructor: Hiram Ring) <i>“Treating text documents like pizza”</i>
SP25	<b>Instructor of Record</b> , DS-GA 1015: Text as Data, New York University
2024	<b>Invited Lecture</b> , DS-GA 3001: Computational Linguistics and Cognitive Science, New York University (Instructor: Tal Linzen) <i>“Transformer-based language model surprisal predicts human reading times best with about two billion training tokens”</i>
2024	<b>Invited Lecture</b> , Metro Early College High School <i>“Guessing meaning and breaking it down (with NACLO problems)”</i>
AU20, SP21	<b>Instructor of Record</b> , LING 3801: Codes and Code Breaking, The Ohio State University Overall student evaluation: 4.69/5 (AU20), 4.86/5 (SP21)
SU17	<b>Teaching Assistant</b> , College English 1, Seoul National University

## SERVICE

---

### ORGANIZING

2024–2025	Workshop on Cognitive Modeling and Computational Linguistics
-----------	--

### REVIEWING

2025	ACL Rolling Review, CogSci 2025, <i>Nature Human Behavior</i>
2024	ACL Rolling Review, CogSci 2024, BlackboxNLP, <i>Journal of Memory and Language</i> , <i>Frontiers in Communication</i> , <i>eLife</i> , <i>Nature Computational Science</i>
2023	ACL Rolling Review
2022	ACL Rolling Review
2021	ACL Rolling Review, 11th Workshop on Cognitive Modeling and Computational Linguistics

### DEPARTMENTAL/UNIVERSITY

2024–2025	MS Admissions Committee, Center for Data Science, New York University
2020, 2022	Student Volunteer, North American Computational Linguistics Olympiad (NACLO)

2019–2024	Laboratory/Computing Committee, Department of Linguistics, The Ohio State University
2019–2020	Treasurer, Student Linguistics Association, The Ohio State University
2018–2019	Speakers Committee, Department of Linguistics, The Ohio State University
2017–2018	Resident Advisor, Gwanak Residence Halls, Seoul National University
2015–2017	International Student Assistant, Gwanak Residence Halls, Seoul National University
OTHER	
2015–2016	K-MOOC Monitoring Team, National Institute for Lifelong Learning, Seoul, Korea
2012–2014	Translator/Interpreter (ROK Air Force), ROK/US Combined Forces Command, Seoul, Korea

## SKILLS

---

Languages	Korean (native), English (bilingual), Japanese (conversational)
Programming Languages	Python, C++, Bash, R, $\text{\LaTeX}$ , Prolog
Frameworks	PyTorch, TensorFlow, HuggingFace, pandas, Matplotlib, Git, Slurm, PBS