

Department of Computer Science



Dr. Hanjun Dai

Research Scientist
Google Research, Brain Team



 **Date: 16 June 2021 (Wednesday)**

 **Time: 10:00am – 11:00am GMT+8 (HKT)** *Zoom details will only be provided to registrants

Scaling up Autoregressive Models for Structured Data

ABSTRACT

Generative modeling remains challenging for large scale discrete structured data like sequences, trees or graphs. The most commonly used model for such data structures is the autoregressive one. However, how to balance between the expressiveness and the computational tractability for the autoregressive models is still a hot research question.

In this talk, we first introduce a scalable autoregressive model BiGG [1] for generating graph structures, where it can reduce the inference cost from $O(n^2)$ to $O(n \log n)$ for a graph with n nodes while still maintaining the full autoregressive capability. Then in the second part we show how this idea could be applied for the self-attention module in Transformers to reduce its asymptotic computation and memory cost, especially for long sequences, while still maintaining the full attention in the unidirectional language modeling scenario. This variant named Combiner [2] can achieve state-of-the-art performance in tasks including ImageNet generative modeling.

References:

- [1] Scalable Deep Generative Modeling for Sparse Graphs, ICML 2020;
- [2] Combiner: Full Attention Transformer with Sparse Computation Cost, in submission;



BIOGRAPHY

Hanjun Dai is currently a research scientist at Google Research, Brain Team. He obtained his PhD degree from Georgia Institute of Technology, advised by Prof. Le Song. His research focuses on deep learning with structured data, including combinatorial optimization, generative modeling, and the corresponding applications. He has published over 30 papers in top-tier conferences and journals, while his work has been recognized by AISTATS 2016 best student paper and best workshop papers in Recsys 2016 and NIPS 2017.

ENQUIRY