

# Byung Soo Jeon

DISTRIBUTED INFERENCE @ NVIDIA | CS PHD @ CMU | U.S. LAWFUL PERMANENT RESIDENT (LPR)

✉ 412-628-3003 | ✉ soojeonml@gmail.com | 🏠 madfunmaker.github.io

## Profile

I am a research-minded engineer specializing in high-performance distributed ML systems for LLMs, with full-stack experience across both training and inference spanning ML compiler optimization, multi-GPU parallelism, operator fusion, and GPU kernel development.

## Professional Experience

Sep 2024 -	<b>Senior System SW Engineer</b> , Multi-GPU inference and performance tech lead for TensorRT compiler	NVIDIA
	<ul style="list-style-type: none"><li>Led design and development of multi-GPU support for TensorRT's compiler backend from scratch</li><li>Optimized context/tensor parallelism end-to-end, from graph-level op fusion to custom CUDA kernels</li></ul>	
Jul - Sep 2024	<b>Staff ML System Engineer</b> , Optimized batching for a distributed LLM inference engine	OctoAI (now NVIDIA)
2017 - 2024	<b>Research Assistant</b> , Thesis: Automated & Portable Machine Learning System	CMU
	<ul style="list-style-type: none"><li>Built ML compiler/distributed runtime that delivers portable, optimized performance for Transformers — up to 1.6x speedup (training) and 1.4x (inference) over baseline ML frameworks (see publications)</li><li>Automated graph-level optimizations, reducing manual tuning effort from days to minutes</li></ul>	
Summer 2023	<b>Research Intern</b> , Benchmarked parallelisms for LLM inference and its implication on HW/SW co-design	Google
Summer 2020	<b>Applied Scientist Intern</b> , Researched on efficient meta-reinforcement learning and active exploration	Amazon
Summer 2019	<b>Applied Scientist Intern</b> , Developed a simulator / RL agent that solves online binpacking problem	Amazon
Jan - Jun 2017	<b>Research Intern</b> , Developed end-to-end multi-modal neural network for music emotion recognition	Naver
2015 - 2016	<b>Researcher</b> , Developed a distributed system / algorithm for billion-scale tensor algebra	KAIST / SNU
2013	<b>Co-founder &amp; SW Engineer</b> , Developed client and server for multiplayer racing mobile game	Funpresso, Inc

## Education

### CMU (Carnegie Mellon University)

PH.D. IN COMPUTER SCIENCE

Pittsburgh, PA

May 2024

### KAIST (Korea Advanced Institute of Science and Technology)

B.S. IN COMPUTER SCIENCE (SUMMA CUM LAUDE)

Daejeon, Korea

Aug 2015

## Publications

### AUTOMATED AND PORTABLE ML SYSTEM

#### Cache Parallelism: Comparative Analysis of Parallelisms in Distributed LLM Inference for Long Sequence Thesis Chapter

Byungsoo Jeon, TIANQI CHEN, ZHIHAO JIA

#### GraphPipe: Improving the Performance and Scalability of DNN Training with Graph Pipeline Parallelism

ASPLOS 2025

Byungsoo Jeon\*, MENGDI WU\*, SUNGHYUN KIM\*, SHIYI CAO\*, SUNGHYUN PARK, NEERAJ AGGARWAL, COLIN UNGER, DAAYAN ARFEEN, PEIYUAN

LIAO, XUPENG MIAO, MOHAMMAD ALIZADEH, GREGORY R. GANGER, TIANQI CHEN, ZHIHAO JIA

#### Collage: Seamless Integration of Deep Learning Backends with Automatic Placement

PACT 2022

Byungsoo Jeon\*, SUNGHYUN PARK\*, PEIYUAN LIAO, SHENG XU, TIANQI CHEN, ZHIHAO JIA

- Integrated to Apache TVM Open-source Project (v0.9.0) | Presented in GTC 2022

#### SRTuner: Effective Compiler Optimization Customization By Exposing Synergistic Relations

CGO 2022

SUNGHYUN PARK, SALAR LATIFI, YONGJUN PARK, ARMAND BEHROOZI, Byungsoo Jeon, SCOTT MAHLKE

### APPLIED ML / RL

#### FactoredRL: Leveraging factored graphs for deep reinforcement learning

NeurIPS 2020

BHARATHAN BALAJI\*, PETROS CHRISTODOULOU\*, XIAOYU LU\*, Byungsoo Jeon, JORDAN BELL-MASTERTON

DeepRL Workshop

#### OBP-RL: Exploring Deep Reinforcement Learning Methods for Online Binpacking Problem

AMLC 2020

Byungsoo Jeon, BHARATHAN BALAJI, SAURABH GUPTA, CHUN YE

Amazon ML Conf

## Dropout Prediction over Weeks in MOOCs by Learning Representations of Clicks and Videos

Byungsoo Jeon\*, NAMYONG PARK\*

AAAI 2020

AI4Edu Workshop

## Dropout Prediction over Weeks in MOOCs via Interpretable Multi-Layer Representation Learning

Byungsoo Jeon\*, NAMYONG PARK\*, SEOJIN BANG\*

AAAI 2020

AI4Edu Workshop

## Time-series Insights into the Process of Passing or Failing Online University Courses using Neural-Induced Interpretable Student States

Byungsoo Jeon, EYAL SHAFRAN, LUKE BREITFELLER, JASON LEVIN, CAROLYN P. ROSE

EDM 2019

## Attentive Interaction Model: Modeling Changes in View in Argumentation

YOHAN JO, SHIVANI PODDAR, Byungsoo Jeon, QINLAN SHEN, CAROLYN P. ROSE, GRAHAM NEUBIG

NAACL 2018

## Music Emotion Recognition via End-to-End Multimodal Neural Networks

Byungsoo Jeon, CHANJU KIM, ADRIAN KIM, DONGWON KIM, JANGYEON PARK, JUNGWOO HA

RecSys 2017

Poster

## DISTRIBUTED SYSTEM / ALGORITHM FOR TENSOR ALGEBRA

### BIGtensor: Mining Billion-Scale Tensor Made Easy

NAMYONG PARK\*, Byungsoo Jeon\*, JUNGWOO LEE, U KANG

CIKM 2016

Demo paper

### SCouT: Scalable Coupled Matrix-Tensor Factorization - Algorithm and Discoveries

Byungsoo Jeon, INAH JEON, LEE SAEL, U KANG

ICDE 2016

### TeGViz: Distributed Tera-Scale Graph Generation and Visualization

ByungSoo Jeon, INAH JEON, U KANG

ICDM 2015

Demo paper

## Teaching

---

Spring 2021 **Machine Learning Systems**, TA (Instructor: Tianqi Chen)

CMU

Spring 2020 **Deep Reinforcement Learning and Control**, TA (Instructor: Katerina Fragkiadaki)

CMU

Spring 2019 **Machine Learning (PhD)**, TA (Instructors: Leila Wehbe, Aaditya Ramdas)

CMU

## Fellowship

---

2022 **Qualcomm Innovation Fellowship**, One of 19 winners in US (\$100k for an year)

Qualcomm

2017 - 2021 **Kwanjeong Scholarship**, One of ~ 50 nationwide outstanding PhD students in STEM (\$30k per year)

KEF

## Skills

---

**ML / Distributed System** C++, Python | TensorRT, PyTorch, TVM, vLLM | CUDA, NCCL, SLURM, MPI