

Selected Publications/Events 2 Video on YouTube: <u>Build Your Own ChatBot with Neural Chat | Intel Software</u> (Oct 2023) Blog published on Medium: Layer-wise Low-bit Weight Only Quantization on a Laptop (Oct 2023) Blog published on Medium: Intel-Optimized Llama.CPP in Intel Extension for Transformers (Oct 2023) Blog published on Medium: Reduce the Carbon Footprint of Large Language Models (Oct 2023) • Blog published on Medium: Empower Applications with Optimized LLMs: Performance, Cost, and Beyond (Sep 2023) • Blog published on Medium: NeuralChat: Simplifying Supervised Instruction Fine-tuning and Reinforcement Aligning for Chatbots (Sep 2023) Intel Innovation'23 Keynote: <u>Intel Innovation 2023 Keynote by Greg Lavender</u> (Sep 2023) • Blog published on Medium: NeuralChat: A Customizable Chatbot Framework (Sep 2023) View Full Publication List. Additional Content 2 Release Information Contribution Guidelines Legal Information Security Policy Apache License **Acknowledgements** 2 • Excellent open-source projects: bitsandbytes, FastChat, fastRAG, ggml, gptq, llama.cpp, lm-evauationharness, peft, trl, streaminglim and many others. • Thanks to all the contributors. Collaborations 2 Welcome to raise any interesting ideas on model compression techniques and LLM-based chatbot development!

Feel free to reach us, and we look forward to our collaborations on Intel Extension for Transformers!

Privacy Contact GitHub Security Status Docs Pricing Training Blog