SW/HW Innovations in Emerging DL Training Systems

Urmish Thakker Engineering Manger, NLP Group

SambaNova

Goldilocks Zone

Too Hot









Trend of SOTA Models



TinyBERT: Distilling BERT for Natural Language Understanding

Xiaoqi Jiao¹*, Yichun Yin²*, Lifeng Shang^{2‡}, Xin Jiang² Xiao Chen², Linlin Li³, Fang Wang^{1‡} and Qun Liu² ¹Key Laboratory of Information Storage System, Huazhong University of Science and Technology, Wuhan National Laboratory for Optoelectronics ²Huawei Noah's Ark Lab

{yin
{che
 DistilBERT, a distilled version of BERT: smaller,
 faster, cheaper and lighter

Victor SANH, Lysandre DEBUT, Julien CHAUMOND, Thomas WOLF Hugging Face {victor,lysandre,julien,thomas}@huggingface.co

Bigger Models





160

180

140

Emerging Hardware and Systems







GRAFHCORE

MYTHIC





Our Mission

Shaping the next-generation ML / DL computing system to accelerate the full model spectrum





How do we break out of the Godilocks Zone?

Fundamental advances required at all layers of the SW/HW stack.



The SambaNova Systems Advantage



Application innovations

High model accuracy

High compute efficiency



Part 1.

Enabling higher compute efficiency

Architecture: Reconfigurable Dataflow Unit (RDU)





Spatial Dataflow Within an RDU

The old way: kernel-by-kernel





SambaFlow eliminates overhead and maximizes utilization



Copyright © 2022 SambaNova Systems, Inc. All rights reserved

Rapid Dataflow Compilation to RDU



SambaFlow Produces Highly Optimized Spatial Mappings





Uncompromised Programmability and Efficiency Breaking out of the programmability vs. efficiency tradeoff curve





The SambaNova Systems Advantage

Achieve low time-to-accuracy



High model accuracy



Part 2. High model accuracy:

+ Pure 16-bit FPU training

Low Precision (< 32-bit) Training

Binarized Neural Networks: Training Neural Networks with Weights and Activations Constrained to +1 or -1

Matthieu Courbariaux*¹ Itay Hubara*² Daniel Soudry³ Ran El-Yaniv² Yoshua Bengio^{1,4} ¹Université de Montréal ²Technion - Israel Institute of Technology ³Columbia University ⁴CIFAR Senior Fellow *Indicates equal contribution. Ordering determined by coin flip. MATTHIEU.COURBARIAUX @ GMAIL.COM ITAYHUBARA @ GMAIL.COM DANIEL.SOUDRY @ GMAIL.COM RANI @ CS.TECHNION.AC.IL YOSHUA.UMONTREAL @ GMAIL.COM

Recurrent Neural Networks With Limited Numerical Precision

Joachim Ott*, Zhouhan Lin[‡], Ying Zhang[‡], Shih-Chii Liu*, Yoshua Bengio^{‡†} *Institute of Neuroinformatics, University of Zurich and ETH Zurich ottj@ethz.ch, shih@ini.ethz.ch [‡]Département d'informatique et de recherche opérationnelle, Université de Montréal [†]CIFAR Senior Fellow {zhouhan.lin, ying.zhang}@umontreal.ca

Training Deep Neural Networks with 8-bit Floating Point Numbers

Naigang Wang, Jungwook Choi, Daniel Brand, Chia-Yu Chen and Kailash Gopalakrishnan IBM T. J. Watson Research Center Yorktown Heights, NY 10598, USA {nwang, choij, danbrand, cchen, kailash}@us.ibm.com

Higher system efficiency, minimal impact on acc. for specific models



Efficiency of Low Precision Floating-point-units (16 vs. 32-bit)



1.5X lower chip area

3X higher energy efficiency

1.5X higher throughput

1. Horowitz. ISSCC 2014

2. Galal et. al. ISCA 2013



Mixed Precision for Generic DL Training (16 + 32 bits FPU)

NVIDIA / apex lines 52.5k		O PyTorch
A PyTorch Extension: Tools for easy mixed precision		Table of Contents
4.7k stars 양 632 forks	$\equiv \uparrow $ TensorFlow	
	TensorFlow Core	AUTOMATIC MIXED PRECISION PACKAGE - TORCH.CUDA.AMP
	TensorFlow > Learn > TensorFlow Cor	re > Guide

Mixed precision

Illusion: 16-bit FPU alone is not enough to maximize model acc.



Can we support only 16-bit FPU on accelerators

&

achieve model acc. matching 32-bit training?

Pure 16-bit (BFloat16) FPU Training





The Accuracy Challenge



Standard 16-bit FPU training degrades model accuracy



The Devil: Nearest Rounding(NR) for Model Weight Updates





Rounding

Update

The Devil: Nearest Rounding (NR) for Model Weight Updates

Theory sketch for least-squares regression

$$\|\boldsymbol{w}_{t} - \boldsymbol{w}^{*}\| \geq \mathcal{O}\left(\boldsymbol{\epsilon} \cdot \min_{j} |\boldsymbol{w}_{j}^{*}|\right)$$
Optimal solution j-th dim of the optimal solution

Inaccurate weight update fundamentally degrades convergence



Stochastic Rounding to the Rescue



Intuition

The expectation of unbiased estimates is as accurate as weights w/o rounding



Kahan Summation as Alternative Enhancement

Auxiliary 16-bit values to track and correct weight update errors from NR







Pure 16-bit training can match 32-bit training in model acc.





Summary

With support for



Accelerators with only 16-bit compute units can match acc. of 32-bit training



The SambaNova Systems Advantage



Application innovations



Part 3. Model Innovations:

Powered by our architecture and algorithm

Computer Vison Evolution of high-resolution Deep Learning



Low-resolution (e.g. cats)

4k images (e.g. Autonomous driving)

50k x 50k (e.g. astronomy, medical imaging, virus, ...



No Compromise High-Res Segmentation



Training w/o information loss from full-image processing



High-Res Pathology with Slide-level Label (TCGA)

Train with Patch label = slide label



Noisy patches limits model accuracy



High-Res Pathology with Slide-level label (TCGA)



16X larger patches \rightarrow 6 Pt higher AUC



Natural Language Processing

Breakthrough efficiency in NLP model online deployment



Distilled tiny Bert model

Short sequence input







Enable up to 11X speedup for online training and inference



Expand Pareto Frontier beyond GPU Higher acc. at higher downstream training and inference throughput (RDU/CPU...)





Scaling Laws for Neural Language Models



Application accuracy improves as the size of the language model increases



Copyright © 2022 SambaNova Systems, Inc. All rights reserved

GPT Family





Copyright © 2022 SambaNova Systems, Inc. All rights reserved

Train Large Language Models, Without Code Changes 1T parameter NLP training with a small footprint and programming ease



https://arxiv.org/pdf/2104.04473.pdf



Enabling Large Model Architectures With a Single System

Order of magnitude performance improvement, an order of magnitude fewer systems



"One Model" 1Trillion Params in a Single System: Same Programming Model



Copyright © 2022 SambaNova Systems, Inc. All rights reserved

