

11. F. Cristian, "Understanding Fault-Tolerant Distributed Systems," *Comm. ACM*, vol. 34, 1993, pp. 56–78.
 12. T.G. Dietterich, "Ensemble Methods in Machine Learning," *Proc. Int'l Workshop Multiple Classifier Systems*, 2000, pp. 1–15.
 13. D. Beyer, *The Future of Machine Intelligence: Perspectives from Leading Practitioners*, O'Reilly, 2016.
 14. R.P. Pothukuchi et al., "Using Multiple Input, Multiple Output Formal Control to Maximize Resource Efficiency in Architectures," *Proc. ACM/IEEE 43rd Ann. Int'l Symp. Computer Architecture*, 2016; doi:10.1109/ISCA.2016.63.
 15. R. Adolf et al., "Fathom: Reference Workloads for Modern Deep Learning Methods," *Proc. IEEE Int'l Symp. Workload Characterization*, 2016; doi:10.1109/IISWC.2016.7581275.
- Yuhao Zhu** is a PhD candidate in the Department of Electrical and Computer Engineering at the University of Texas at Austin. Contact him at yzhu@utexas.edu.
- Vijay Janapa Reddi** is an assistant professor in the Department of Electrical and Computer Engineering at the University of Texas at Austin. Contact him at vj@ece.utexas.edu.

The Design and Evolution of Deep Learning Workloads

ROBERT ADOLF
SAKETH RAMA
BRANDON REAGEN
GU-YEON WEI
DAVID BROOKS
 Harvard University

..... The past decade has witnessed the reemergence of a connectionist approach to solving several classes of challenging artificial intelligence problems. This family of strategies is collectively known as representation learning, hierarchical learning, or, most popularly, deep learning. The success of deep learning, like many facets of cognitive computing, is the result of a confluence of progress in three separate areas, rather than a single, monumental breakthrough. These three areas include the collection and curation of massive datasets, advances in machine learning algorithms, and the ever-increasing power of computational hardware. These three phenomena form a virtuous cycle. Success in one area facilitates growth in the other two, along with increasing demand for it. For instance, the landmark win of a deep neural network at the ImageNet

Large Scale Visual Recognition Challenge in 2012 was the result of a massive new set of training data¹ (two orders of magnitude larger than its closest predecessor), several clever novel modifications to a many-layer convolutional neural network,² and use of high-performance hardware (among the first to leverage GPUs for deep learning).

That cognitive computing should be characterized as much by data and hardware as algorithms is not surprising: the very definition involves learning by example at scale. However, it does suggest that carrying out research in this field is perhaps unique, in that one cannot make ample headway without considering all three aspects. This multidimensional constraint is felt keenly in the creation and curation of representative workloads for deep learning problems. Benchmarks and proxy applications must strike a bal-

ance between simplicity and faithful reproduction, accurately capturing all fundamental aspects of the programs they represent while remaining easy to understand, use, and transform. Doing this across several axes is challenging, even more so given the frenetic pace of innovation and upheaval in the field. We believe the right approach is first to design workloads that capture the unique aspects of deep learning models, data, and implementations, and then to embrace change and plan for continuous evolution.

Design

Building good deep learning benchmarks means getting three things right: choosing the right models, respecting the impact of data, and faithfully reproducing unique implementation details. We

Table 1. The Fathom workloads

Model	Dataset	Style	Purpose and legacy
Seq2Seq	WMT-15	Supervised, recurrent	Direct language-to-language sentence translation. State-of-the-art accuracy with a simple, language-agnostic architecture.
MemNet	bAbI	Supervised, memory network	Facebook's memory-oriented neural system. One of two novel architectures that explore a topology beyond lattices of neurons.
Speech	TIMIT	Supervised, recurrent, fully connected	Baidu's speech-recognition engine. Proved purely deep-learned networks can beat hand-tuned systems.
Autoenc	MNIST	Unsupervised, fully connected	Variational autoencoder. An efficient, generative model for feature learning.
Residual	ImageNet	Supervised, convolutional	Image classifier from MSR Asia. Dramatically increased the depth of convolutional networks. ILSVRC 2015 winner.
VGG	ImageNet	Supervised, convolutional, fully connected	Image classifier demonstrating the power of small convolutional filters. ILSVRC 2014 winner.
AlexNet	ImageNet	Supervised, convolutional, fully connected	Image classifier. Watershed for deep learning by beating hand-tuned image systems at ILSVRC 2012.
DeepQ	Atari ALE	Reinforcement, convolutional, fully connected	Atari-playing neural network from DeepMind. Super-human performance on many Atari2600 games, without any preconceptions.

present our case in the context of our experience in designing Fathom, a set of reference workloads for deep learning (see Table 1).³

Models

The most visible decision for a workload suite is the choice of which models to include. In Fathom, we used three criteria to select eight models from a wide array of candidates: representativeness, diversity, and impact. The first is clear: our choices should reflect the best of what the deep learning community has come up with. Because there are many models that could rightly claim this status, the need to limit the size implies a need for diversity; each model should bring something unique to the table. Finally, "impact" reflects the degree to which a particular technique has changed the landscape of deep learning research. We cannot predict the future of deep learning, so we instead tried to choose methods that have imparted fundamental lessons to the work that came after—lessons that will continue to be relevant even as subsequent research builds on them.

Datasets

Data also plays a central role in machine learning workloads, even for architects

and system designers. Although it is true that some fundamental deep learning techniques (such as matrix math, convolution, and backpropagation) are somewhat agnostic to the values of their inputs, the role of data is broader. Many problem domains are heavily affected by how their data is being used. For instance, most supervised learning problems have two different operational modes: training, which involves massive amounts of fixed data and an emphasis on throughput, and inference, which involves a stream of unseen data and a lean toward latency. The same model can exhibit different computational characteristics depending on which environment it is used in. Additionally, much of the research in executing deep learning problems centers around exploiting features unique to neural networks or a specific model structure. Sparsity in weight values, batchsize-convergence tradeoffs, and the degree of downsampling in pooling operations are just a few features that depend heavily on the characteristics of the inputs under consideration.

Implementation

Writing reference workloads involves a balancing act between faithfully mim-

icking praxis while preserving ease of use for researchers. One example of this is the widespread adoption of high-level programming frameworks such as TensorFlow or Torch. These frameworks provide two main benefits: they abstract the underlying hardware interface away from the programmer, and they provide tested libraries of kernels that act as a productivity multiplier. They have changed the development landscape, largely for the better, and it is no longer possible to create a realistic set of deep learning workloads without taking them into account. All eight Fathom models are written on top of TensorFlow. On the other hand, no such consensus has been reached on the layout of learning models, the staging and preprocessing of data, or the mechanisms that drive high-level control flow. It is common to see two implementations of the same model that are almost unrecognizable. Because these choices are more a matter of taste than any fundamental property of deep learning models, Fathom imposed a standard structure and set of interfaces over all its workloads. This greatly simplified cross-model instrumentation, data collection, and experimentation for its users.

Evolution

Deep learning is a field in flux, and a workload suite designed for such an environment must have a plan for adapting. Graceful evolution is an extension of good design: the core principle is to understand which aspects of a workload are intrinsic to the field and which are a product of the current state of the art. For instance, although it's likely that the set of models included in Fathom will change, their selection criteria will not. One convenient way to understand this idea is to look backwards at the developments leading to the present—that is, to understand what changes Fathom would have had to weather had it been released earlier.

Models

All but one of the current Fathom workloads were published since 2014, but most have predecessors that would have been replaced. For instance, DeepSpeech was a breakthrough in pure deep learning speech recognition, but many prior state-of-the-art systems used a combination of hidden Markov models and neural networks. The more interesting change would have been the introduction of read-write networks. Memory networks and neural Turing machines both arrived in 2014, and while no work is built in a vacuum, it is unlikely Fathom would have had something similar. The same is probably true for reinforcement learning: the concept has a long history, but it needed a champion, DeepMind, to make it a core theme in deep learning. Fathom would probably have grown in size over the past several years, in addition to needing to replace its speech model. This is a trend that is almost certain to continue. Even now, it seems likely that additional advances in speech and language processing will require both of Fathom's recurrent models to evolve, and new architectures like binary-valued networks are on the horizon.

Datasets

Surprisingly, most of Fathom's current datasets are relatively stable. ImageNet

has not seen radical changes since its introduction, and MNIST and TIMIT have long histories. The largest change would have been the introduction of the Arcade Learning Environment—the Atari emulator used by Fathom's DeepQ model. Although ALE's inputs are not substantially different from older image datasets, its use and integration are. Training and inference with deep-Q learning is a substantially different beast because it requires two-way, online communication. The underlying trend here is refreshingly optimistic: datasets change because deep learning is improving. While ImageNet will probably remain in Fathom, it seems likely that a new source of visual data will augment it, because recent models have surpassed human performance. Additionally, it seems likely that new datasets using video, graphs, or mixed-mode inputs could merit inclusion.

Implementation

Superficially, an older Fathom would appear substantially different, because TensorFlow was not made public until late 2015. However, the use of high-level frameworks has been a clear trend for several years, so it is likely that Theano, Caffe, or Torch would have been used instead. All four frameworks share similarity in their designs and interfaces. The largest difference would have been the effort required to implement some of the models. Although today's frameworks are all converging on support for most of the techniques Fathom uses, that was not true three years ago. Many of the primitives were implemented in only one library, and most of the analysis tools constructed to characterize the Fathom workloads would have been substantially more difficult. This is largely a result of maturity. Today, deep learning frameworks have larger user bases and more developers, and most common operations are well-supported in all platforms. Given this convergent evolution, it is unlikely that Fathom will need to change its implementation framework in the future. On the other hand, Fathom is fac-

ing a clear need to adapt to another trend: fixed-precision and packed arithmetic. The use of non-floating-point math and limited precision has been known for decades, but most of the work in deep learning has been focused on improving accuracy, discovering new models, and applying them to new problems. As deep learning applications are deployed in mainstream scenarios, however, efficiency and speed have become a central concern. Many frameworks have introduced some form of packed arithmetic, supported by vector instructions on CPUs and more recently by double-speed, half-width operations on GPUs. This trend is only increasing in importance, and Fathom will need to adopt some form of it to keep up.

Deep learning is a protean field, and workloads for it must be living projects. This is a challenge for maintainers as well as researchers, but it also reflects the success of the virtuous cycle that drives it. Evolution implies that all three facets—models, data, and hardware—are still moving forward in lockstep. Moreover, our experience with Fathom suggests that there are consistent principles underpinning the process that can guide the requisite adaptation. We look forward to the bright future of deep learning, and we believe that accurate, practical, and fluid workloads will continue to play an important role in its progress. MICRO

References

1. O. Russakovsky et al., "ImageNet Large Scale Visual Recognition Challenge," *Int'l J. Computer Vision*, vol. 115, no. 3, 2014, pp. 211–252.
2. A. Krizhevsky, I. Sutskever, and G.E. Hinton, "ImageNet Classification with Deep Convolutional Neural Networks," *Advances in Neural Information Processing Systems (NIPS)*, 2012, pp. 1097–1105.
3. R. Adolf et al., "Fathom: Reference Workloads for Modern Deep Learning

Methods," *Proc. IEEE Int'l Symp. Workload Characterization (IISWC)*, 2016, pp. 1–10.

Robert Adolf is a PhD candidate in the School of Engineering and Applied Sciences at Harvard University. Contact him at rdadolf@seas.harvard.edu.

Saketh Rama is a PhD student in the School of Engineering and Applied Sciences at Harvard University. Contact him at rama@seas.harvard.edu.

Brandon Reagen is a PhD candidate in the School of Engineering and Applied Sciences at Harvard University. Contact him at reagen@fas.harvard.edu.

Gu-Yeon Wei is the Wei Gordon McKay Professor of Electrical Engineering and Computer Science at Harvard University. Contact him at guyeon@eecs.harvard.edu.

David Brooks is the Haley Family Professor of Computer Science at Harvard Uni-

versity. He is an editorial board member of *IEEE Micro*. Contact him at dbrooks@eecs.harvard.edu.

myCS

Read your subscriptions through the myCS publications portal at <http://mycs.computer.org>.



Can You Invent a Better World through Technology?

Challenge Accepted

Computer Society Global Student Challenge

Create a solution, based on the IEEE Computer Society 2022 report, that will solve a real-world issue.

Over US\$2,000 in Prizes!
1st place gets US\$1,500 and will be honored at the Annual Awards Banquet in Phoenix, AZ in June 2017

Submission Deadline: 1 April 2017

Enter the challenge at computer.org/studentchallenge