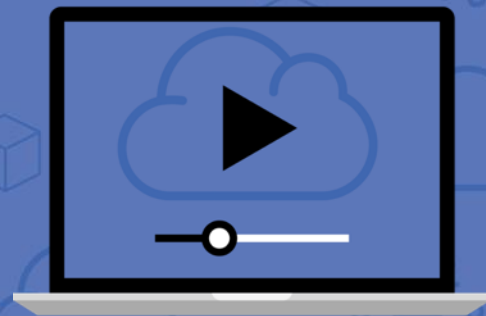




Machine Learning: From Notebook to Production with Amazon Sagemaker

Julien Simon, AI Evangelist, EMEA
@julsimon





Welcome to Amazon.com Books!

*One million titles,
consistently low prices.*

(If you explore just one thing, make it our personal notification service. We think it's very cool!)

SPOTLIGHT! -- AUGUST 16TH

These are the books we love, offered at Amazon.com low prices. The spotlight moves **EVERY** day so please come often.

ONE MILLION TITLES

Search Amazon.com's [million title catalog](#) by author, subject, title, keyword, and more... Or take a look at the [books we recommend](#) in over 20 categories... Check out our [customer reviews](#) and the [award winners](#) from the Hugo and Nebula to the Pulitzer and Nobel... and [bestsellers](#) are 30% off the publishers list...

EYES & EDITORS, A PERSONAL NOTIFICATION SERVICE

Like to know when that book you want comes out in paperback or when your favorite author releases a new title? Eyes, our tireless, automated search agent, will send you mail. Meanwhile, our human editors are busy previewing galleys and reading advance reviews. They can let you know when especially wonderful works are published in particular genres or subject areas. Come in, [meet Eyes](#), and have it all explained.

YOUR ACCOUNT

Check the status of your orders or change the email address and password you have on file with us. Please note that you **do not** need an account to use the store. The first time you place an order, you will be given the opportunity to create an account.

Amazon.com, 1995

« Two Decades of Recommender Systems at Amazon.com » (2017)



Amazon is well-known for personalization and recommendations, which help customers discover items they might otherwise not have found. In this update to our original article, we discuss some of the changes as Amazon has grown.

Two Decades of Recommender Systems at Amazon.com

It's two decades now. Amazon.com has been building a store for every customer. Each person who comes to Amazon.com sees it differently, because it's individually personalized based on their interests. It's as if you walked into a store and the shelves started rearranging themselves, with what you might want moving to the front, and what you're unlikely to be interested in shuffling farther away. From a catalog of hundreds of millions of items, Amazon.com's recommendation pick a small number of items you might enjoy based on your current context and your past behavior. The algorithms aren't magic; they simply chase with you what other people have already discovered. The algorithm does all the work. It's computer helping people help other people, implicitly and anonymously. Amazon.com launched item-based collaborative filtering in 1998, enabling recommendations at a previously unseen scale for millions of customers and a catalog of millions of items. Since we wrote about the algorithm in IEEE Internet Computing in 2003, it has seen widespread use across the web, including YouTube, Netflix, and many others. The algorithm's success has been due to its simplicity, scalability, and other surprising and useful

recommendations, as well as desirable properties such as updating immediately based on new information about a customer and being able to explain why it recommended something in a way that's easily understandable. What we described in our 2003 IEEE Internet Computing article has found many challenges and seen much development over the years. Here, we describe some of the updates, improvements, and adaptations for item-based collaborative filtering, and offer our view on what the future holds for collaborative filtering, recommender systems, and personalization.

The Algorithm As we described in 2003, the item-based collaborative filtering algorithm is straightforward. In the mid-1990s, collaborative filtering was generally user-based, meaning the first step of the algorithm was to search across other users to find people with similar interests (such as similar purchase patterns), then look at what items those similar users found that you haven't found yet. In short, our algorithm begins by finding related items for each item in the catalog. The term "related" could have several meanings here, but at this point,

Two Decades of Recommender Systems at Amazon.com

Standing the Test of Time

A part of recognizing IEEE Internet Computing for its 20 years in publication, I remembered to the editorial board that we pick one of our magazine articles that, over the past 20 years, has withstood the "test of time." In selecting an article, we evaluated the time to more than 20 candidate articles that reported on "evergreen" research areas over the past two decades and then assessed these articles based on downloads from IEEE Xplore, citations, and mentions of the work in popular press. This information was presented to a committee consisting of previous Editors in Chief for the magazine. I would like to thank the selection committee from the editorial board — led by Aron Jorgensen, including Fred Douglass, Robert Hertz, Michael Hahn, Charles Perkins, Michael Kabanovich, and Harshvardh Singh. This committee deliberated on the top three articles by evaluating each work's previous exposure within the context of its sustained importance in the future.

It's my pleasure to recognize the committee's official "Test of Time" winner, an industry article titled "Amazon.com Recommendations: Item-to-Item Collaborative Filtering" by Greg Linden.

Brian Smith, and Jeremy Turk, from the January/February 2003 issue of IC, doi:10.1109/IC.2003.1207464. Reprints were after the publication of this article. It shows 102 downloads from IEEE Xplore to date, with more than 12,754 downloads since January 2015. The article currently shows 628 citations in Google Scholar. To align with the selection committee recommendation as an industry article, as I align with the magazine's focus of accessibility to academic, research, and education populations. In addition to recognizing this article, we asked the authors to create this retrospective piece discussing research and insights that have resounded since publishing their winning "Test of Time" article, while preparing into the future. Going forward, the magazine hopes to celebrate a "Test of Time" article every 2-3 years. I hope that you enjoy this retrospective article, and please take a moment to congratulate Greg Linden, Brian Smith, and Jeremy Turk.

—H. Brian Babble
Editor-in-Chief, IEEE Internet Computing
Professor and Distinguished Professor, Central University

let's loosely define it as "people who buy one item are unusually likely to buy the other." So, for every item i , we want every item j , that was purchased with unusually high frequency by people who bought i .

Since this related items table is built, we can generate recommendations quickly as a series of lookups. For each item that's part of this customer's current context and previous interests, we look up the related items, combine them to yield the most likely items of interest, filter out items already seen or purchased, and then we are left with the items to recommend.

This algorithm has many advantages over the older user-based collaborative filtering. Most importantly, the stability of the computation is more robust to a batch build of the related items — and the computation of the recommendations can be done in real time as a series of lookups. The recommendations are high quality and useful, especially given enough data, and remain competitive in perceived quality even with the newer algorithms created over the last two decades. The algorithm scales to hundreds of millions of users and tens of millions of items without sampling or other techniques that reduce the quality of the recommendations. The algorithm updates immediately on new information about a person's interests. Finally, the recommendations can be explained in an

iterative way as arising from a list of items the customer remembers purchasing.

In 2003, Amazon.com, Netflix, YouTube, and More

By the time we published in IEEE in 2003, item-based collaborative filtering was widely deployed across Amazon.com. The homepage prominently featured recommendations based on your past purchases and items bought in the store. Search result pages recommended items related to your search. The shopping cart recommended other items to add to your cart, perhaps inspired by

its handle in the store, or perhaps comments in what you were already considering. Most importantly, the stability of the computation is more robust to a batch build of the related items — and the computation of the recommendations can be done in real time as a series of lookups. The recommendations are high quality and useful, especially given enough data, and remain competitive in perceived quality even with the newer algorithms created over the last two decades. The algorithm scales to hundreds of millions of users and tens of millions of items without sampling or other techniques that reduce the quality of the recommendations. The algorithm updates immediately on new information about a person's interests. Finally, the recommendations can be explained in an

G.D. Linden, J.A. Jacobi, and E.A. Benson, Collaborative Recommendations Using Item-to-Item Similarity Mappings, US Patent 6,266,649, to Amazon.com, Patent and Trademark Office, 2001 (filed 1998).

<https://www.computer.org/csdl/mags/ic/2017/03/mic2017030012.html>

amazonrobotics



amazon echo





INTRODUCING
amazon go



Democratization of AI

APPLICATION SERVICES



Amazon
Rekognition



Amazon Polly



Amazon Lex

Amazon Comprehend



Amazon
Rekognition Video

Amazon Transcribe

Amazon Translate

PLATFORM SERVICES

Amazon SageMaker

AWS DeepLens

Amazon EMR

FRAMEWORKS AND INTERFACES

Deep Learning
AMI

Apache MXNet

Caffe2

CNTK

PyTorch

TensorFlow

Theano

Torch

Keras

Gluon

More AI/ML is built on AWS than anywhere else

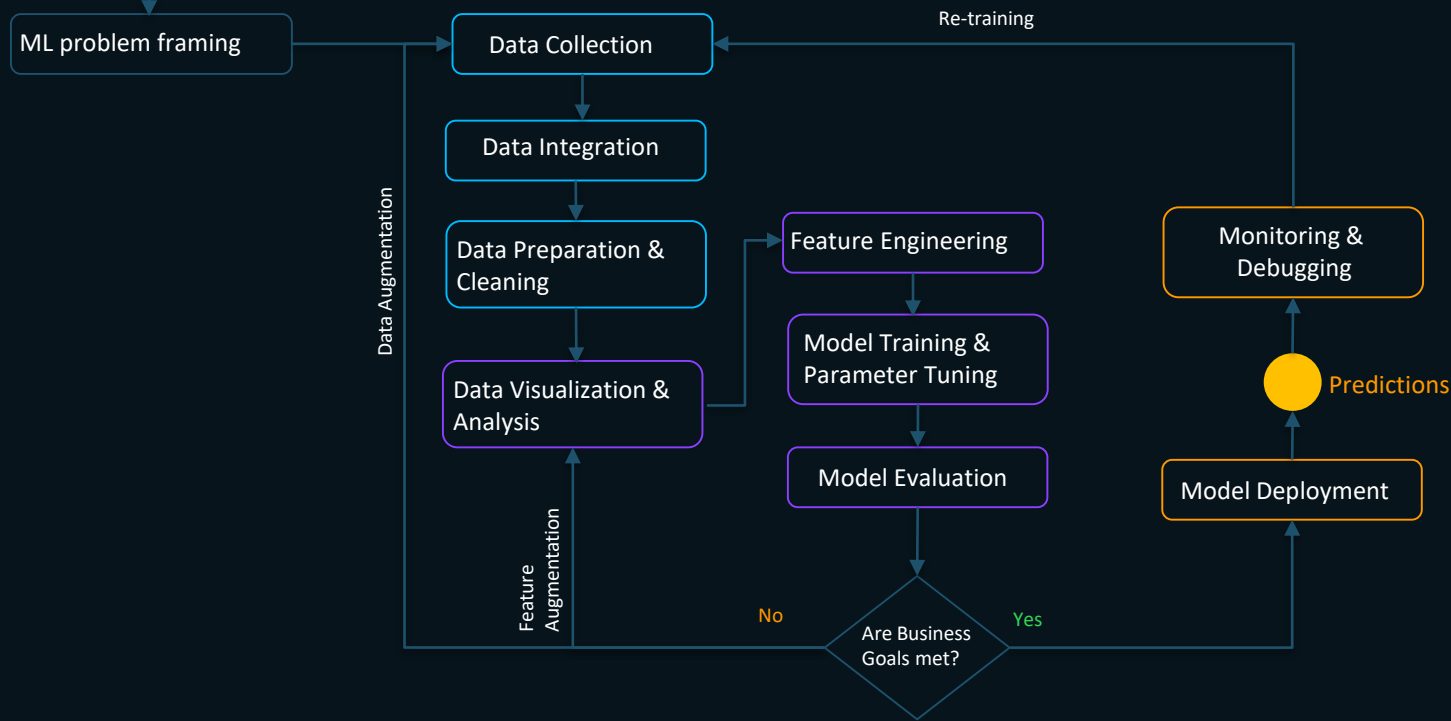


More AI/ML is built on AWS than anywhere else



Business Problem –

The Machine Learning Process

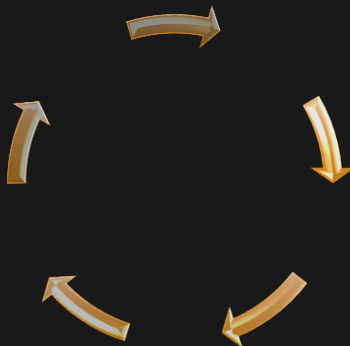


Amazon SageMaker

Build

Pre-built notebook
instances

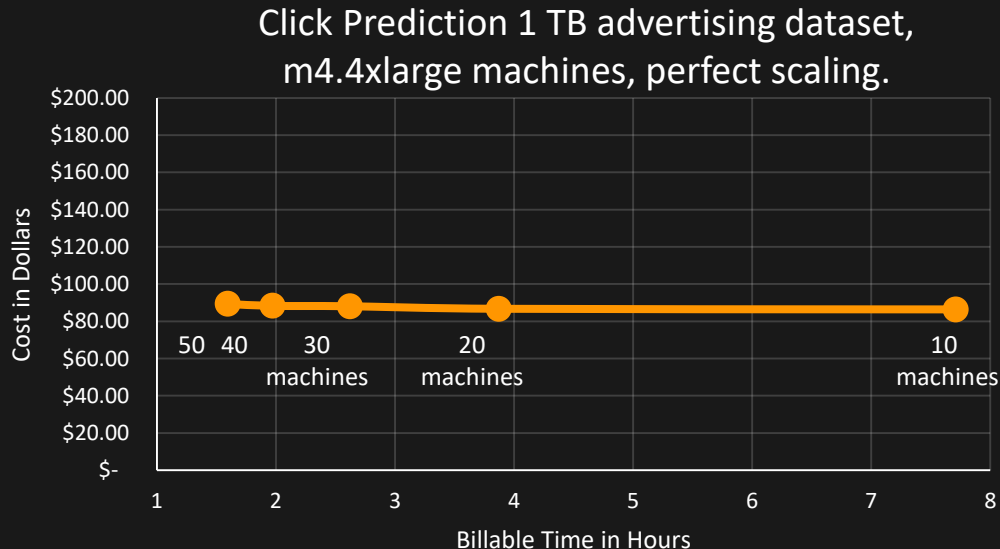
Highly-optimized
machine learning
algorithms



Factorization Machines

$$\tilde{y} = w_0 + \langle w_1, x \rangle + \sum_{i,j>i} x_i x_j \cdot \langle v_i, v_j \rangle$$

	Log_loss	F1 Score	Seconds
SageMaker	0.494	0.277	820
Other (10 Iter)	0.516	0.190	650
Other (20 Iter)	0.507	0.254	1300
Other (50 Iter)	0.481	0.313	3250



Spectral LDA

The New York Times

U.S.

High-Tech Industry, Long Shy of Politics, Is Now Belle of Ball

By LIZETTE ALVAREZ DEC. 26, 1999

Correction Appended

At a time when Congress is bitterly divided and unable to reach consensus on issues like gun control and health care, Democrats and Republicans are happily reaching across party lines to pass legislation backed by high-tech companies.

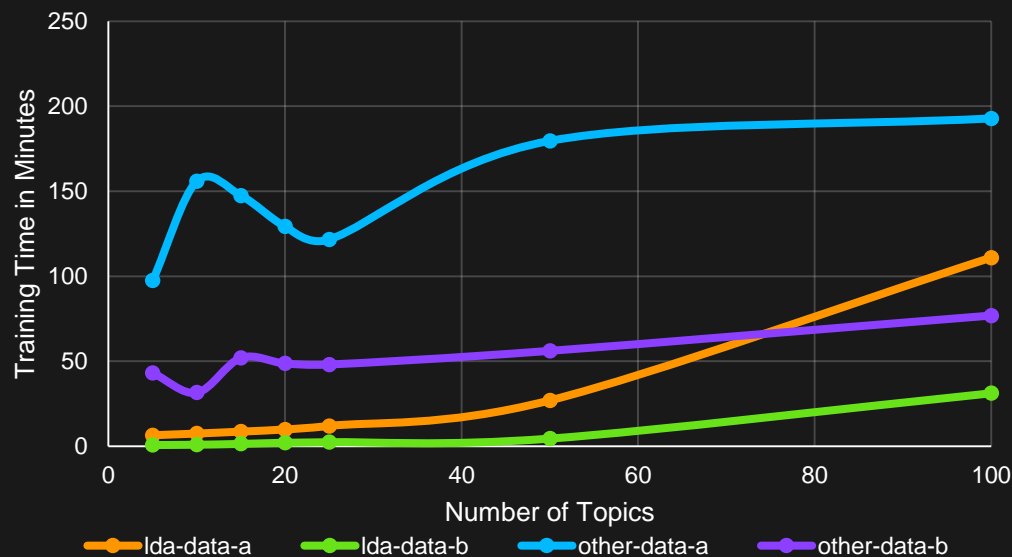
The high-tech industry, at the same moment, is lavishing new attention on Washington and changing its once aloof posture toward the federal government.

Republicans and Democrats are both eager to win the loyalties of high-tech companies and executives, knowing that they represent untold jobs, wealth and ultimately votes and campaign contributions.

For its part, the industry has realized that the federal government can do its members as much harm as good. Microsoft, and its battle with the Justice Department, along with a spate of other threatened legal problems, drilled this point home.

"Microsoft was a poster child for our industry," said Connie Correll, director of communications for the Information Technology Industry Council, a trade organization that represents America Online, Dell and I.B.M., among others.

Training Time vs. Number of Topics



Sequence to Sequence

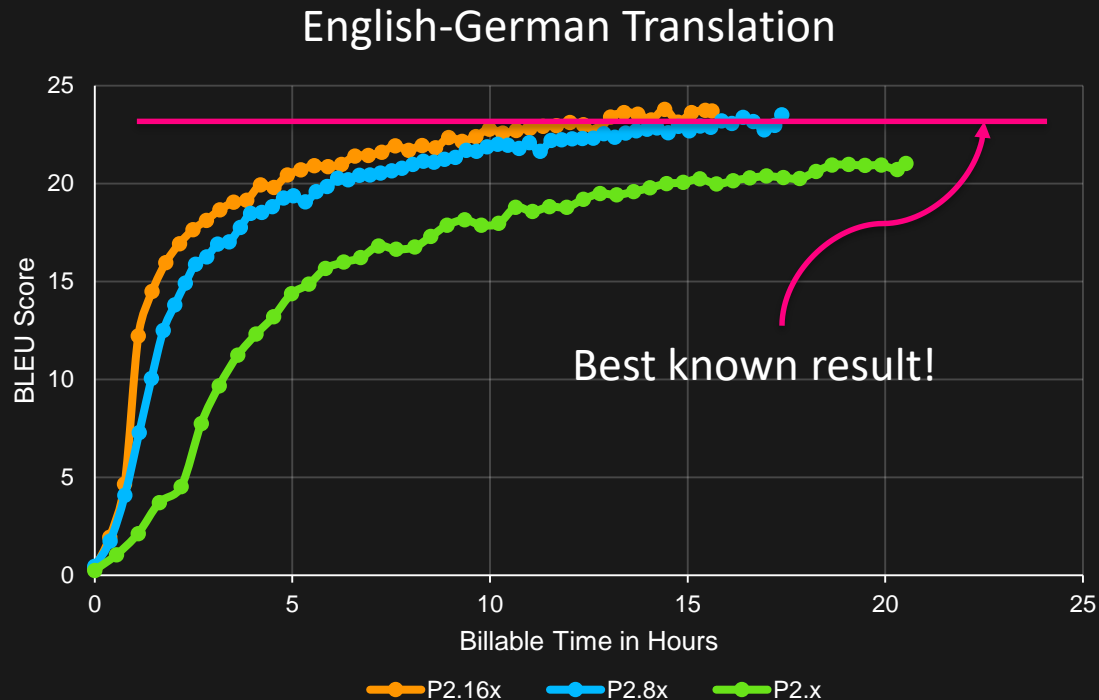
Based on Sockeye and
Apache MXNet.

Supports multi-GPU training.

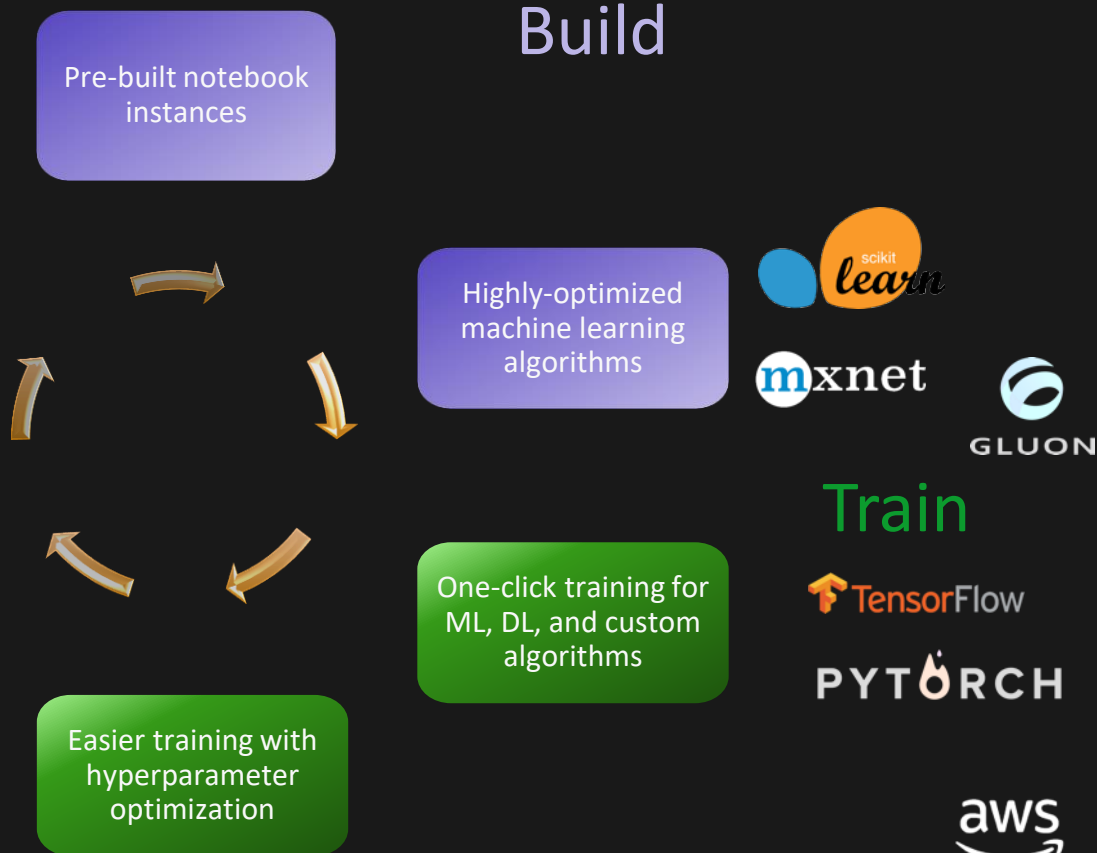
Can be used for Neural Machine
Translation.

Supports both RNN/CNN
as encoder/decoder

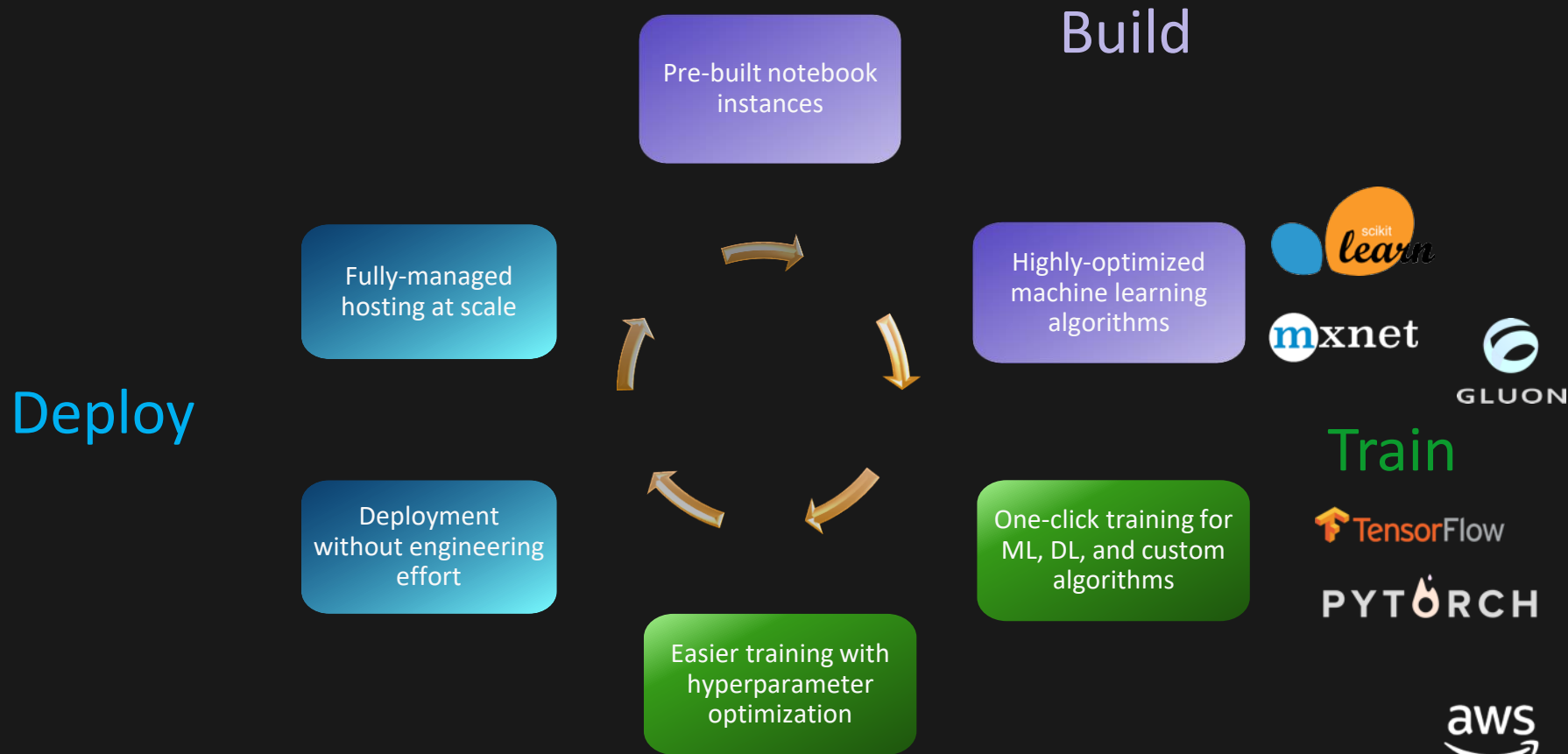
<https://arxiv.org/abs/1712.05690>
<https://github.com/aws-labs/sockeye>



Amazon SageMaker



Amazon SageMaker

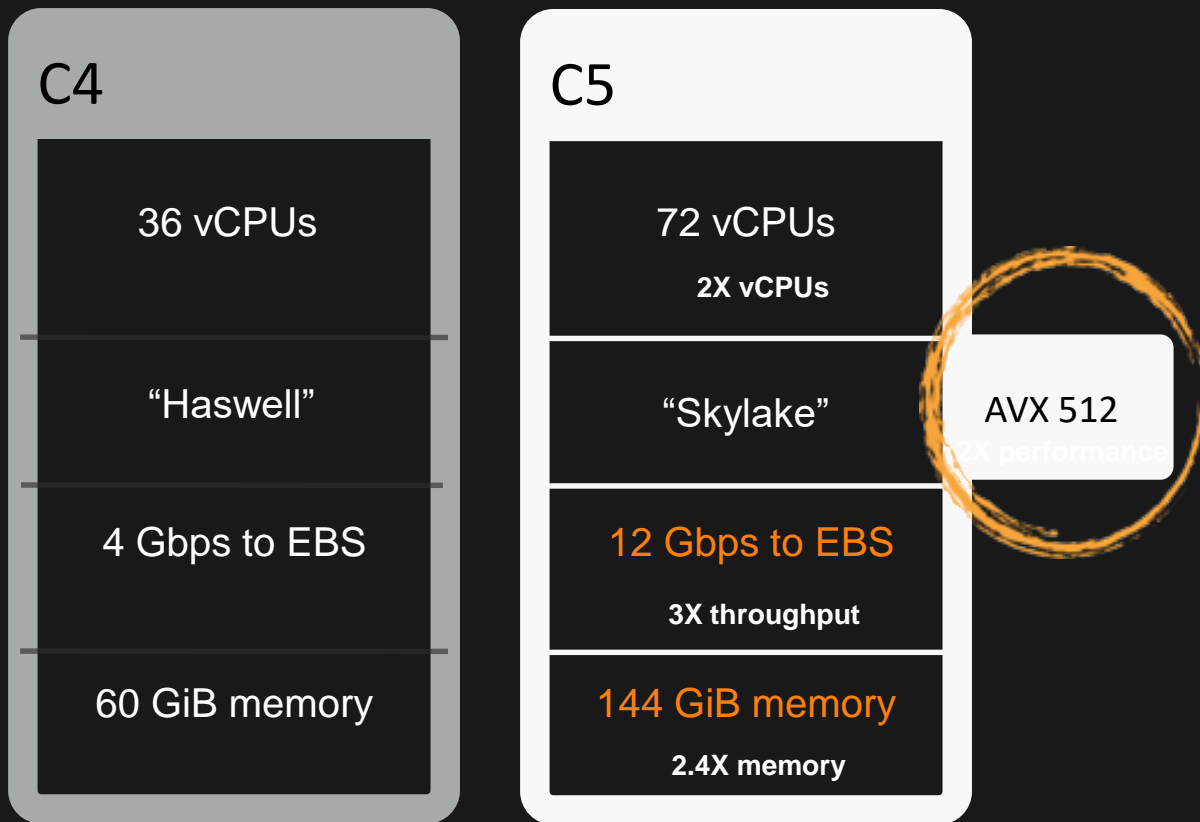


Infrastructure – Amazon EC2 P3 Instances

The fastest, most powerful GPU instances in the cloud

- Up to eight NVIDIA Tesla V100 GPUs
- 1 PetaFLOPs of computational performance – *14x better than P2*
- 300 GB/s GPU-to-GPU communication (NVLink) – *9X better than P2*
- 16GB GPU memory with 900 GB/sec peak GPU memory bandwidth

Infrastructure – Amazon EC2 C5



DEMO

Resources

<https://aws.amazon.com/machine-learning>

<https://aws.amazon.com/blogs/ai>

<https://aws.amazon.com/sagemaker>

<https://github.com/aws/sagemaker-python-sdk>

<https://github.com/aws/sagemaker-spark>

<https://github.com/aws-labs/amazon-sagemaker-examples>

An overview of Amazon SageMaker

<https://www.youtube.com/watch?v=ym7NEYEx9x4>

<https://medium.com/@julsimon>



Thank you!

Julien Simon, AI Evangelist, EMEA
@julsimon

