intel® **xeon**®

# Using 4th Gen Intel® Xeon® Processors, Taboola Improves AI-Driven Content Recommendation Engines

## The company's 12,000 servers and ten data centers help deliver tens of billions of personalized content recommendations every day.

### Intel® Technology Highlights

- Optimizing on 4th Gen Intel® Xeon® processors reduces response time and improves model accuracy

- Intel® Advanced Matrix Extensions improve AI capabilities on the CPU

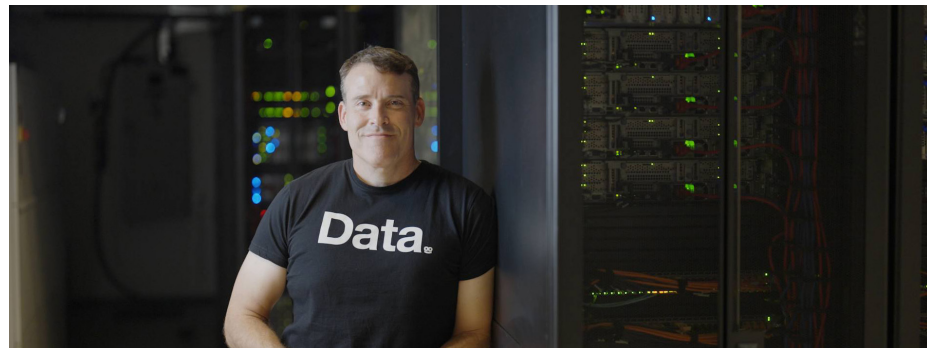- Intel® AVX-512 enables more requests per second per core

intel. **XEON**

**Taboola**

With the average internet user now spending close to 7 hours per day online[1] consumers browsing their favorite sites are becoming more and more resistant to traditionally targeted pop-ups and banner ads. It's critical for advertisers to be able to target this content as accurately as possible to keep even jaded visitors engaged, without slowing down the user experience.

That's why many companies are turning to artificial intelligence (AI)-driven content recommendation engines to determine the most appropriate content to show. Machine leaning algorithms use a plethora of contextual data, including recently viewed content, trending topics, and even location and time of day, to predict what online consumers are most likely to find interesting on an individual basis.

Taboola, a global leader in this space since 2007, uses an AI-based prediction engine at the edge to deliver targeted content and recommendations most relevant to each visitors' unique preferences. With 12,000 servers and ten data centers, Taboola's solutions connect over a billion customers worldwide every month with advertisers such as Honda and Adidas. The result: tens of billions of personalized content recommendations every day, seamlessly integrated into the pages of publishers such as USA Today or MSN.

Advertisers and retailers rely on partners like Taboola to instantly deliver targeted, relevant, and visually appealing recommendations to compliment consumers' mobile experience. "We help advertisers turn regular users into paying customers by providing a service that gives a richer kind of experience from the internet," says Ariel Pisetzky, Taboola's VP of Information Technology & Cyber.

The heart of Taboola's solution is a neural network based on the open source TensorFlow Serving (TFS) framework.



Ariel Pisetzky, Taboola's VP of Information Technology & Cyber, in one of Taboola's ten data centers.

Architected on top of TensorFlow, TFS employs a client-server workflow to deliver recommendations. When a TFS server gets a request from a client, it runs the client data through a pretrained Taboola neural network model and returns the result.

Taboola wanted to increase application throughput without sacrificing efficiency, so they turned to Intel and the 4th Gen Intel® Xeon® processors.

## Challenge

Taboola provides an average of 10 recommendations per page on approximately 4 billion webpages a day. "These numbers really add up," Pisetzky says. "To do all of this, to create and manufacture these recommendations, we need to be super-efficient—not just in terms of costs, but also the impact on the environment" in terms of energy usage.

"We are the cloud for our customers, and we need a CPU that is flexible, versatile, and efficient, so it can perform multiple jobs," Pisetzky says, "It needs to be able to scale up and down through the day in terms of power use and have enough memory to do in-memory calculations that we need done really fast."

As a relatively small company, Taboola needs to balance power and costs by thoughtfully utilizing their technology stack to optimize efficiency and performance across the lifetime of their equipment. Extracting maximum performance per core is key.

"Unlike larger companies, we don't have the luxury to engineer our own solutions, to create our own chips or write our own operating system," Pisetzky says. "We want to spend our research and development cost around optimizing for both advertiser and publisher success."

To serve more content recommendations, Taboola wanted to increase throughput on their AI-driven recommendation engine. Throughput is crucial to scale the ability to match users with the brand and editorial content that's most interesting and relevant to them. Ideally this would happen while keeping client latency below 100 ms.

## Solution

Intel worked with Taboola to optimize and benchmark their prediction algorithm on 4th Gen Intel® Xeon® processors. The result was higher overall performance: average maximum throughput on Intel Xeon 8480 processors was 1.74x higher than on Intel Xeon 8380 processors.[2]

Part of the performance boost came from the Intel® Advanced Matrix Extensions (Intel® AMX) accelerator, which improves AI capabilities on the CPU. This makes it ideal for workloads that include image recognition, recommendation systems, and natural-language processing. Intel also leveraged higher core counts and memory bandwidth on the Intel Xeon 8480 processor and applied the latest TFS with Intel optimizations.

The Intel® AVX-512 accelerator, which increases the size of a CPU's register, also speeded up the process, Pisetzky says. "With Intel AVX-512, I can run more requests per second per core per CPU, and I can load the servers more." Even though competitors may offer more cores per chip, without the Intel AVX-512 instruction set, they need a higher core count to reach the same level of performance. "So you have to look not only at the core count, but what these cores can do."

"It's great to have a hardware partner like Intel with the know-how and engineering capacity to provide a host of solutions, and the knowledge of how to best utilize this hardware for many years to come," Pisetzky says. Optimizing on the latest-generation Intel platforms has helped Taboola reduce response time and improve model accuracy, as well as minimize operational costs as server efficiency improves, he adds.

"Our business model only works when the publisher is successful and when the advertiser is successful," Pisetzky says. "We want our customers to succeed, and working with Intel really helps me base my business around my customers."

Get more information about Taboola
Learn more about 4th Gen Intel® Xeon® processors