

## Make AI models faster, smaller, and less expensive to run.

ByteShape enables organizations to deploy more capable AI on existing hardware while reducing latency, memory requirements, energy consumption, and infrastructure cost.

<b>LOWER COST</b> <b>Improve AI economics</b> Reduce infrastructure, memory, and energy requirements.	<b>MORE CAPABILITY</b> <b>Use existing hardware</b> Deploy larger or higher-quality models on available systems.	<b>MORE CONTROL</b> <b>Own the deployment</b> Run in cloud, on-premises, at the edge, or in sovereign infrastructure.
---	--	---

### AUTOMATED MODEL + HARDWARE OPTIMIZATION

ByteShape automatically optimizes numerical precision, value representation, and execution code for each target model, workload, and hardware platform.

Our system explores fine-grained datatype and quantization choices across tensors and value groups, balancing model quality against customer-defined deployment objectives.

<b>QUALITY</b> Fidelity	<b>LATENCY</b> Response time	<b>THROUGHPUT</b> e.g., Tokens / sec
<b>MEMORY</b> Model footprint	<b>ENERGY</b> Power use	<b>COST</b> Deployment economics

### HOW IT WORKS

- 01 Specify**  
Define the model, target hardware, quality threshold, and deployment objective.
- 02 Optimize**  
ByteShape searches the design space and measures quality and performance.
- 03 Deploy**  
Receive an optimized model and execution path that fit your workflow.

### WHY BYTESHAPE

<b>AUTOMATED SEARCH</b> <b>Beyond fixed recipes and manual tuning</b> Systematically learns a broader optimization space using measured quality and performance.	<b>OBJECTIVE-DRIVEN</b> <b>Built around your deployment target</b> Optimizes against the quality, latency, throughput, memory, energy, and cost targets that matter.
<b>BROAD COMPATIBILITY</b> <b>Models, formats, workloads, and devices</b> Supports integer, floating-point, and microscaling representations across diverse architectures (not only LLMs) and targets.	<b>FUTURE-READY</b> <b>Designed to adapt as AI changes</b> Operates at the level of tensors, datatypes, data movement, and compute kernels - not one fixed model family.

### BUILT FOR REPEATABLE DEPLOYMENT ACROSS THE AI STACK

<b>AI MODEL COMPANIES</b> Ship optimized releases across more devices and price points.	<b>PLATFORMS &amp; INFRASTRUCTURE</b> Improve utilization and offer stronger performance per dollar.	<b>ENTERPRISES &amp; PUBLIC SECTOR</b> Deploy capable AI while retaining control of systems and data.
--	---	--

### PROVEN TECHNICAL FOUNDATION

Peer-reviewed research at MLSys, demonstrated from edge devices to high-end GPUs.  
Commercial model: annual enterprise subscriptions and optimization engagements.

### START A CONVERSATION

Bring us a model, target hardware, and deployment objective.  
[contact@byteshape.com](mailto:contact@byteshape.com)