

Aurora González Vidal
University of Murcia

Alexander Isenko
Technical University of Munich

K.R. Jayaram
IBM Thomas J. Watson Research Center



On Serving Image Classification Models

Contact: aurora.gonzalez2@um.es



This study forms part of the ThinkInAzul programme and was supported by MCIN with funding from European Union NextGenerationEU (PRTR-C17.11) and by Comunidad Autónoma de la Región de Murcia – Fundación Séneca



- 1 Introduction
- 2 Background
- 3 Methodology
- 4 Results
- 5 Future Work

Introduction

Background

Methodology

Results

Future Work

1 Introduction

2 Background

3 Methodology

4 Results

5 Future Work

Introduction

Background

Methodology

Results

Future Work

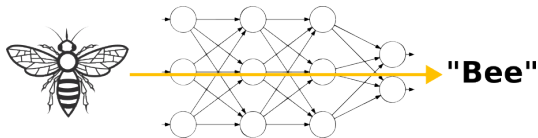
In deep learning applications, up to 90 percent of the infrastructure cost for developing and running an ML application is spent on inference.

Needs: scalable, guarantee high system goodput, and maximize resource utilization.

Intention: Set the foundations for model inference serving in serverless computing environments

Objective: analyse the factors independently and together to build up a generalizable optimization model to assist in scheduling decisions

Use case: Image classification inference because its many applications such as e-commerce and retail (Amazon or Pinterest), social media such as instagram, autonomous vehicles, medical image analysis etc



Introduction

Background

Methodology

Results

Future Work

1 Introduction

2 **Background**

3 Methodology

4 Results

5 Future Work

Types of inference according to deadline guarantees.

- “Hard” Real-time Inference
- “Soft” Real-time
- Relaxed Inference
- Best-effort Inference

Equipment: TPU, GPU, CPU, etc.

Our study case: 1 GPU (NVIDIA A100 with 40 GB of VRAM),
“Soft” Real-time and Relaxed Inference.

1 Introduction

2 Background

3 Methodology

4 Results

5 Future Work

- Selection of an image classification model: EfficientNet-B0
- Creation of dummy images with different input sizes
- Measuring inference times (repeated) over the different input sizes and mini-batch sizes looking for dependencies (for later on defining functions)
- Hardware monitoring ¹ (164 features including network bandwidth, disk read/write bandwidth and counters, CPU parameters, memory utilization, GPU (pynvml and torch): temperature, memory fragmentation, etc.)
- Proposition of mathematical models for the optimization of the inference process

¹<https://github.com/circuit/py-hardware-monitor>

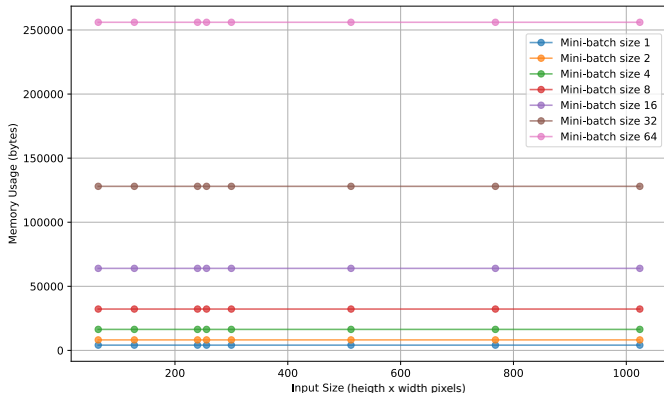
1 Introduction

2 Background

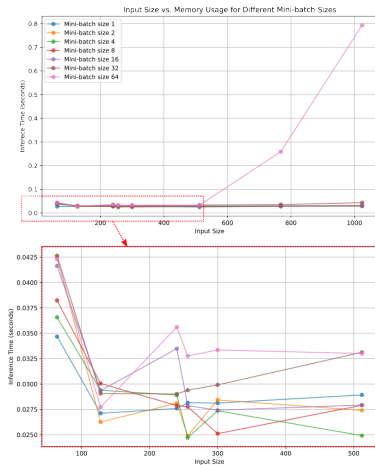
3 Methodology

4 Results

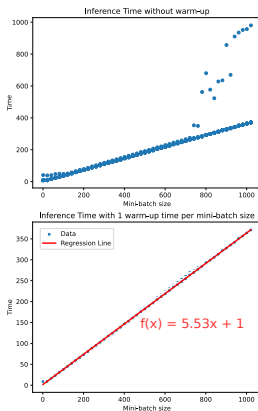
5 Future Work



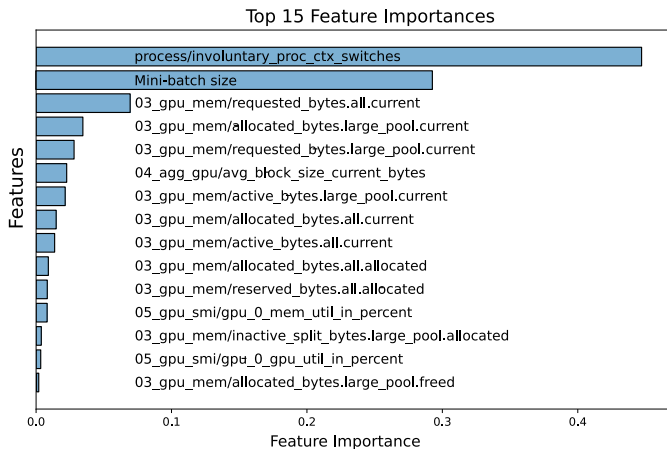
Memory usage using different image input sizes and mini-batch sizes



Memory usage using different image input sizes and mini-batch sizes



Inference time using different mini-batch sizes without considering warm-up (above) and considering warm-up (below) with fixed input size = 224



15 most important features to determining first inference time / warm up

Optimization definitions

Decision variables:

- t_i : The number of times GPU_{*i*} is used (an integer).
- mbs_i : The mini-batch size chosen for GPU_{*i*} (an integer).
- N_G : The number of GPUs to be used (an integer)

The constants:

- T : The total available time. This should not be exceeded by any of the GPUs, given that they work in parallel (a decimal number).
- N : The number of images that need to be processed in total in the given time (an integer).
- $NGPU$: The maximum number of GPUs available (an integer)
- M_i : The maximum number of times GPU_{*i*} can be used (a constant)
- $Size_i$: The images' input size for GPU_{*i*}

The functions:

- L_i : Latency per mbs_i for GPU_{*i*}
- W_i : Warm-up time for GPU_{*i*}
- MB_i : The maximum mini-batch size for GPU_{*i*} (a function of $Size_i$).

$$\begin{aligned}
 \text{mín} \quad & N_G \\
 \text{s.t.} \quad & \text{Maximum}_i (W_i(\text{mbs}_i) + t_i \cdot L_i(\text{mbs}_i)) \leq T \\
 & \sum_i (t_i + 1) \cdot \text{mbs}_i \geq N \\
 & 1 \leq \text{mbs}_i \leq MB_i \quad \text{for all } i \\
 & 0 \leq t_i \leq M_i \quad \text{for all } i \\
 & 1 \leq N_G \leq N_{GPU}
 \end{aligned} \tag{1}$$

$$\begin{aligned}
 \text{máx} \quad & NGPU \times \sum_i (t_i + 1) \cdot \text{mbs}_i \\
 \text{s.t.} \quad & \text{Maximum}_i (W_i(\text{mbs}_i) + t_i \cdot L_i(\text{mbs}_i)) \leq T \quad (2) \\
 & 1 \leq \text{mbs}_i \leq MB_i \quad \forall i \\
 & 0 \leq t_i \leq M_i \quad \forall i
 \end{aligned}$$

- 1 Introduction
- 2 Background
- 3 Methodology
- 4 Results
- 5 Future Work

Future Work:

- Optimal Mini-Batch Determination
- Resource Management and Load Times
- Concurrency and Cost-Energy Limits
- Versatility and Heterogeneous Serving
- Resolution of the optimization models
- Adaptation and Integration

Aurora González Vidal

University of Murcia

Alexander Isenko

Technical University of Munich

K.R. Jayaram

IBM Thomas J. Watson Research Center



On Serving Image Classification Models

Contact: aurora.gonzalez2@um.es



Funded by
the European Union



This study forms part of the ThinkInAzul programme and was supported by MCIN with funding from European Union NextGenerationEU (PRTR-C17.11) and by Comunidad Autónoma de la Región de Murcia – Fundación Séneca

