

Motivation

Baur and Strassen's result, 1983

The arithmetic complexity of evaluating a rational function's derivative is at most 5 times the complexity of function evaluation.

It has now been fifteen years of extensive and global empirical DNN training with nonsmooth components. It was founded on two assumptions:

- backpropagation outputs a gradient almost.
- the process is fast.

Motivation: extends the Baur-Strassen's result to the nonsmooth case.

Automatic differentiation in Machine learning

Given a training set $\{(x_i, y_i)\}_{i=1\dots N}$, the supervised training of a neural network f consists in minimizing the empirical risk:

$$\min_{\theta \in \mathbb{R}^P} J(\theta) := \frac{1}{N} \sum_{i=1}^N \ell(f(x_i, \theta), y_i) \quad (1)$$

where $\theta \in \mathbb{R}^P$ are the network's weight parameters and ℓ is a loss function. In general, f is a composition of nonsmooth functions.

- **Automatic differentiation (AD):** A program that evaluates derivatives of numeric functions expressed as computer programs in an efficient and accurate way.

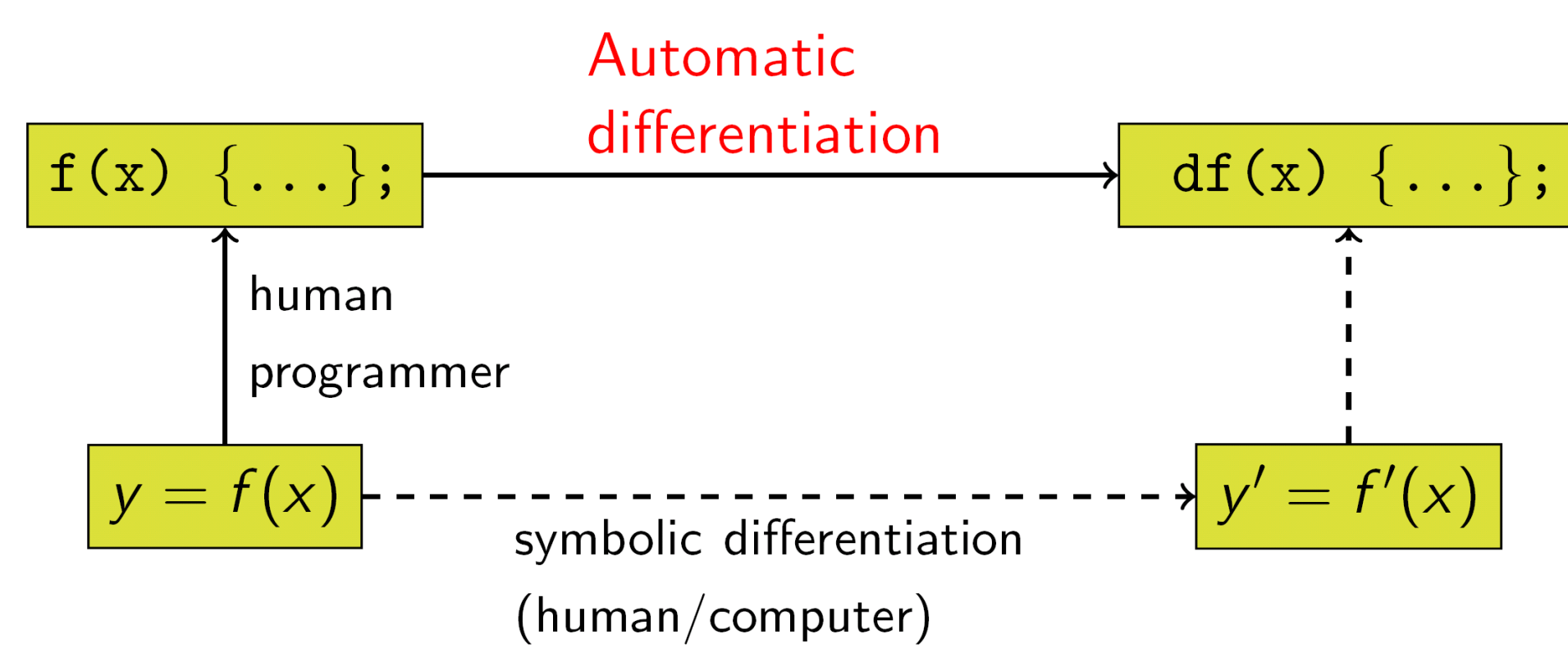


Figure: How automatic differentiation relates to symbolic differentiation

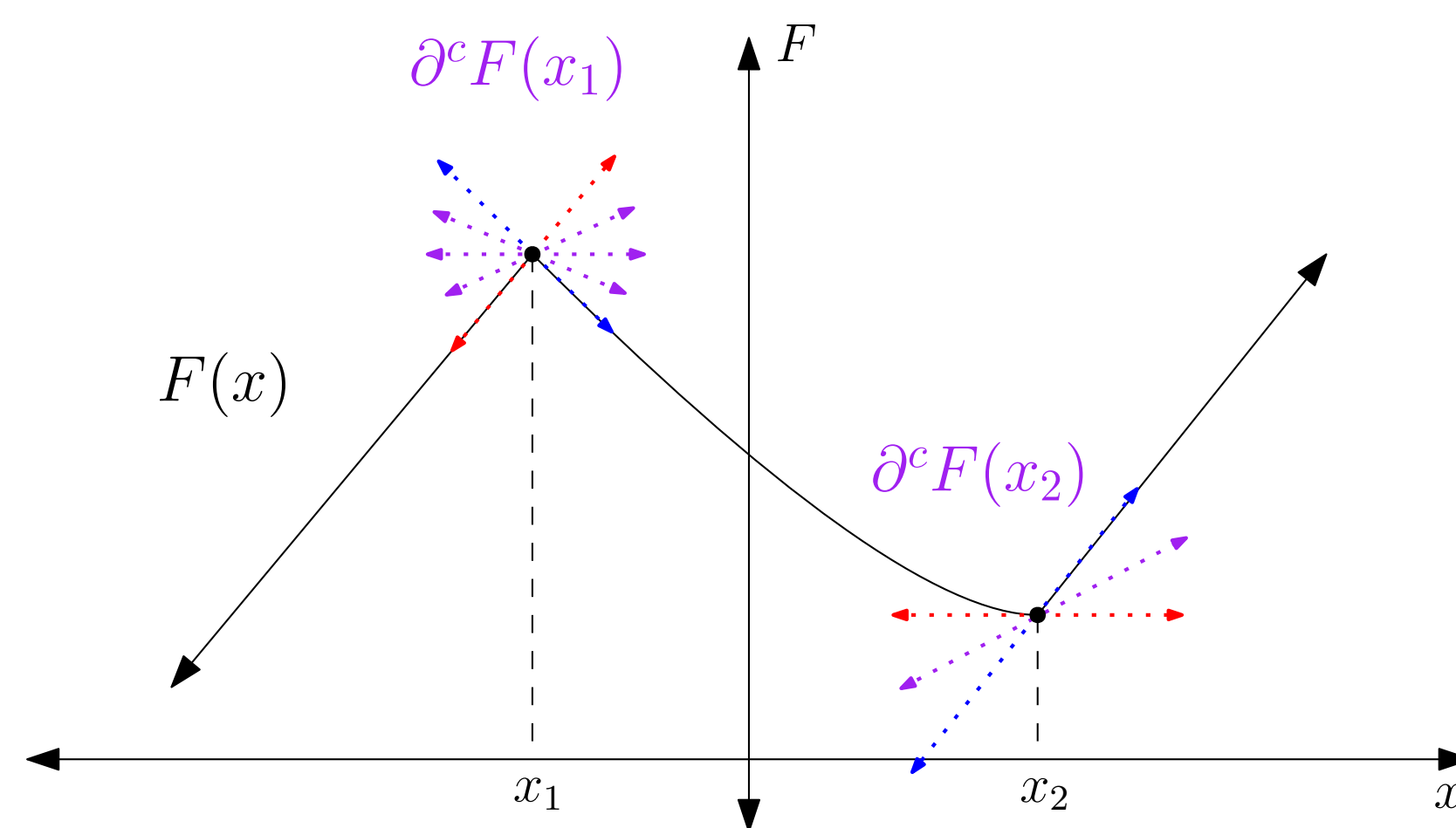
To solve (1), we should use **AD** to compute gradients (in the smooth case) or surrogate gradients (in the nonsmooth case).

Clarke gradients : a nonsmooth oracle

Given a locally Lipschitz continuous function $F : \mathbb{R}^p \rightarrow \mathbb{R}$, the *Clarke subdifferential* of F is

$$\partial^c F(x) = \text{conv} \left\{ \lim_{k \rightarrow +\infty} \nabla F(x_k) : x_k \in \text{diff}_F, x_k \xrightarrow{k \rightarrow +\infty} x \right\} \quad (2)$$

where diff_F is the full measure set where F is differentiable and ∇F is the standard gradient.



Notations

Let $F : \mathbb{R}^p \mapsto \mathbb{R}$ be a locally Lipschitz function.

- \mathcal{D} : collection of elementary operations used to compute F .
- \mathcal{D}' : collection of elementary operations used to compute F and a surrogate gradient of F .
- P : program which computes F using operations on \mathcal{D} .
- $\text{backprop}(P)$: program that computes F and its backpropagation.
- $\text{cost}(P)$: number of \mathcal{D} operations required by the program P .
- $\text{cost}(\text{backprop}(P))$: number of \mathcal{D}' operations required by the program $\text{backprop}(P)$.

Nonsmooth AD with conservative gradients

How does backprop works ?

Consider a locally Lipschitz function $F : \mathbb{R}^p \rightarrow \mathbb{R}$ with a m -compositional representation implemented by a program P

$$F = g_1 \circ \dots \circ g_m.$$

- For each i, x , choose $d_i(x) \in \partial^c g_i(x)$.
- Ex : $g_i = \text{ReLU}$ and take $d_i(0) = 0$ (Tensorflow, Pytorch)

Chain-rule the d_i 's:

$$d_F(x) := d_1(g_2(\dots(g_m(x))\dots)) \times d_2(g_3(\dots(g_m(x))\dots)) \dots \times d_m(x)$$

$\Rightarrow \text{backprop}(P)$ computes $d_F(x)$.

This is how PyTorch and TensorFlow work.

The chain-rule, which is required for AD, usually fails for Clarke subgradients.

\Rightarrow **introduce the conservative gradients.**

Cheap conservative gradient

Let $F : \mathbb{R}^p \mapsto \mathbb{R}$ be a locally Lipschitz function and P a program who compute F using a dictionary \mathcal{D} composed by path differentiable operations. If $F = g_1 \circ \dots \circ g_m$ and each g_i are operations on \mathcal{D} , then :

- F is path differentiable,
- $\text{cost}(\text{backprop}(P)) \leq \omega_b \times \text{cost}(P)$ where ω_b is a constant.

Computational properties of conservative gradients vs others nonsmooth AD frameworks

Others alternative AD approaches

Computational overhead ratio

Minimum value of the quotient of the cost required to evaluate a program and "its" derived program by the cost to evaluate merely the program.

Let $F : \mathbb{R}^p \rightarrow \mathbb{R}$ be a locally Lipschitz function.

Others alternative implementable AD approaches:

- Try to evaluate elements of $\partial^c F$, based on directional derivatives - Khan and Barton (2012;2013;2015)
- Successive local approximations of F , based on lexicographic derivatives
- Computing Clarke subgradients using forward AD

Problem: All these procedures either require to evaluate p directional derivatives.

ReLU networks

Given a set of matrices $M_1 \in \mathbb{R}^{p_1 \times p}$, $M_2 \in \mathbb{R}^{p_2 \times p_1}$, \dots , $M_{L-1} \in \mathbb{R}^{p_{L-1} \times p_{L-2}}$, $M_L \in \mathbb{R}^{1 \times p_{L-1}}$ we consider the associated ReLU network $F : \mathbb{R}^p \rightarrow \mathbb{R}$:

$$F : x \mapsto M_L \text{ReLU}(M_{L-1} \text{ReLU}(\dots M_1 x)).$$

Link between p directional derivatives and matrix multiplication

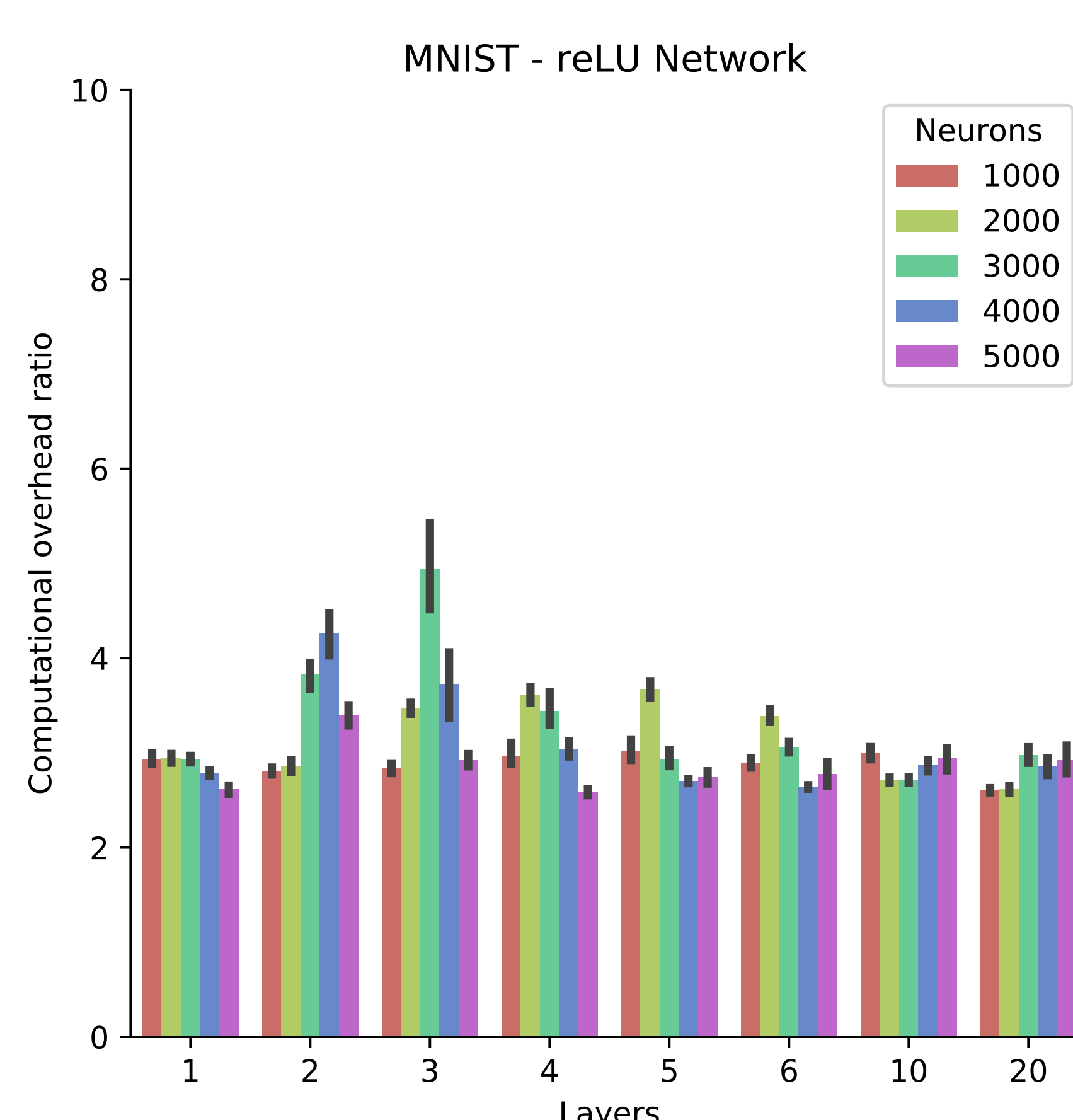
Let $F : \mathbb{R}^p \rightarrow \mathbb{R}$ be a ReLU network function, computational cost defined over \mathbb{R} by circuit complexity ("number of operations"):

- ① $c(p) := \text{cost}(p \times p \text{ matrix multiplication})$
- ② $c(p) \geq p^2$
- ③ $\text{cost}(p \text{ directional derivatives of } F) \geq c(p)$

\Rightarrow **suffers from computational overhead scaling linearly in p**

Applications with ReLU networks

Computational overhead ratio of MLP Cross Entropy with MNIST.



Computational overhead ratio of MLP Cross Entropy with MNIST according to the number of layers/neurons and the batch size.

