# GlobalDataLoader in Multi DeepLearning Task

Xie Jian

I2EC, ICS, NJU

March 9, 2021

# Table of Contents

**1** Introduction

**2** Global DataLoader
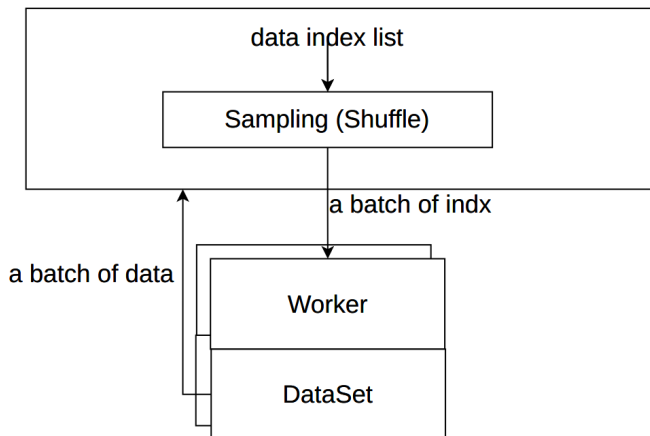
**3** Experiment

# Table of Contents

# DataLoader in Pytorch

# Problem: Repeated Reading and Processing

## Situation

To compare the performance of different algorithms, Many DeepLearning tasks are training in the same Dataset.

## Problem

Every task has its own DataLoader. So the data will be repeatedly read and processed by different tasks.

## Result

As the number of tasks increases, so does the training time. And what increases is the time to load the data
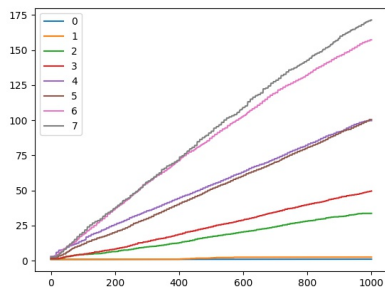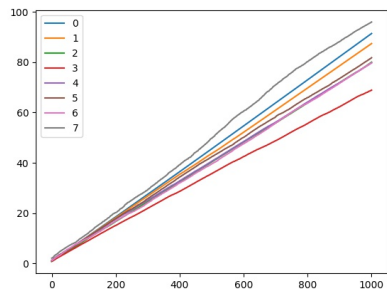
# Experiment

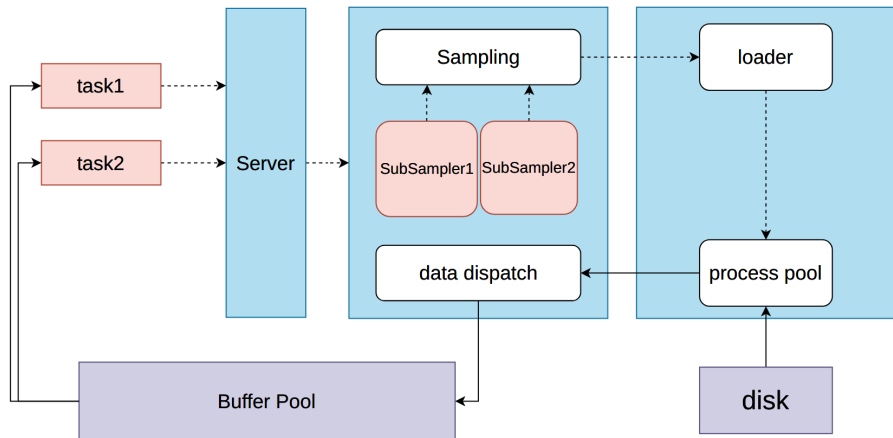

Figure: data loading time



Figure: data training time

# Table of Contents

**1** Introduction

**2** Global DataLoader

**3** Experiment

# Architecture

# Sampling: problem description

### Defination

For a single task, the sampler needs to select an element from the index set $S$.

Similarly, for multiple tasks, the sampler needs to select some elements $\{s_1, s_2, ...\}$ from multiple sets $\{S_1, S_2, ...\}$

### Requirements

- The index in the set $S$ should be randomly sampled. $p(s_i) = \frac{1}{|S_i|}$

- Duplicate indexes need to be merged. Maximize $p(s_i = s_j)$

- There can be no problem of "starvation"

# Independently Sampling Algorithm

## Assumption

There are two sets: $S_1, S_2$, and their length is $n_1, n_2$

The intersection set of them is $S_i$, whose length is $n_i$

We divide the set $S_1$ into $S_i$ and $S_{d1} = S_1 - S_i$

We divide the set $S_2$ into $S_i$ and $S_{d2} = S_2 - S_i$

## Example

$S_1 = \{1, 2, 3, 4, 5\} = \{1, 2, 3\} \cup \{4, 5\}$

$S_2 = \{1, 2, 3, 6, 7\} = \{1, 2, 3\} \cup \{6, 7\}$

# Independently Sampling Algorithm

## Algorithm

- $S_1$:
  - $step_{11}$: randomly select a set from $S_i$ and $S_{d1}$
  - $step_{12}$: if the set is $S_{d1}$, randomly select a element1 from $S_{d1}$
  - $step_{13}$: if the set is $S_i$, randomly select a element1 from $S_i$
- $S_2$:
  - $step_{21}$: randomly select a set from $S_i$ and $S_{d2}$
  - $step_{22}$: if the set is $S_{d2}$, randomly select a element2 from $S_{d2}$
  - $step_{23}$: if the set is $S_i$, randomly select a element2 from $S_i$

## Probability

$$p(element1) = \frac{1}{n_1}, p(element2) = \frac{1}{n_2}$$
$$p(element1 = element2) = \frac{n_i}{n_1 * n_2}$$

$$(1)$$

# Dependently Sampling Algorithm I

## Idea

The $step_{13}$ is same as $step_{23}$. We can merge them.

## Algorithm

- $S_1$:
    - $step_{11}$: randomly select a set from $S_i$ and $S_{d1}$
    - $step_{12}$: if the set is $S_{d1}$, randomly select a element1 from $S_{d1}$
    - $step_{13}$: if the set is $S_i$, go to $step_{33}$
- $S_2$:
    - $step_{21}$: randomly select a set from $S_i$ and $S_{d2}$
    - $step_{22}$: if the set is $S_{d2}$, randomly select a element2 from $S_{d2}$
    - $step_{23}$: if the set is $S_i$, go to $step_{33}$
- $step_{33}$: randomly select a element from $S_i$

# Dependently Sampling Algorithm I

**Probability**

$$p(element1) = \frac{1}{n_1}$$

$$p(element2) = \frac{1}{n_2} \qquad (2)$$

$$p(element1 = element2) = \frac{n_i}{n_1} * \frac{n_i}{n_1} = \frac{n_i^2}{n_1 * n_2}$$

**Example**

$S_1 = \{1, 2, 3\} \cup \{4, 5\}; \ S_2 = \{1, 2, 3\} \cup \{6, 7\}$

- As for $S_1$, select set $\{1, 2, 3\}$ with probability 0.6
- As for $S_2$, select set $\{1, 2, 3\}$ with probability 0.6
- Finally, randomly sampling element in $\{1, 2, 3\}$
- $p(element1 = element2) = 0.6 * 0.6 = 0.36$

# Dependently Sampling Algorithm I

## Idea

The $step_{11}$ and $step_{21}$ are similar. We can merge them.

## Algorithm

- $step_1$: randomly select a set from $S_i$ and $S_{d1}$
- $S_1$:
    - $step_{12}$: if the set is $S_{d1}$, randomly select a element1 from $S_{d1}$
    - $step_{13}$: if the set is $S_i$, go to $step_{33}$
- $S_2$:
    - $step_{22}$: if the set is not $S_d1$, randomly select a element2 from $S_{d2}$
    - $step_{23}$: if the set is $S_i$, go to $step_{33}$
- $step_{33}$: randomly select a element from $S_i$

# Problem

As for $S_1$:

$$p_1(S_i) = \frac{n_i}{n_1}$$

$$p_1(S_{d1}) = 1 - \frac{n_i}{n_1}$$

(3)

As for $S_2$:

$$p_2(S_i) = \frac{n_i}{n_2}$$

$$p_2(S_{d1}) = 1 - \frac{n_i}{n_2}$$

(4)

if $n_1 \neq n_2$,
then $p_1(S_i) \neq p_2(S_i)$ and $p_1(S_{d1}) \neq p_2(S_{d2})$

# Case1: n1 < n2

## Problem

$$(p_1(S_i) = \frac{n_i}{n_1}) > (p_2 S_i = \frac{n_i}{n_2}) \tag{5}$$

## Approach

So in $step_1$, when select $S_i$, it should be changed $S_{d2}$ in probability of $p$. The equation is

$$p_2(S_i) = p_1(S_i) * (1 - p) = \frac{n_i}{n_2}$$
$$p_2(S_{d2}) = p_2(S_{d1}) + p_1(S_i) * p = 1 - \frac{n_i}{n_2} \tag{6}$$

then

$$p = 1 - \frac{n_1}{n_2}$$

# Dependently Sampling Algorithm II

## Algorithm

- $step_1$: randomly select a set $S$ from $S_i$ and $S_{d1}$
- $S_1$:
    - $step_{12}$: if the $S$ is $S_{d1}$, randomly select a element1 from $S_{d1}$
    - $step_{13}$: if the $S$ is $S_i$, go to $step_3$
- $S_2$:
    - $step_{22}$: if the $S$ is $S_{d2}$, randomly select a element2 from $S_{d2}$
    - $step_{23}$: if the $S$ is $S_i$
        - randomly let $S = S_{d1}$ in probability of $\frac{n_1}{n_2}$
        - if $S$ is $S_{d2}$, randomly select a element2 from $S_{d2}$
        - if $S$ is $S_i$, go to $step_3$
- $step_3$: randomly select a element from $S_i$

# Dependently Sampling Algorithm II

## Probability

If in $step_1$, the selected set $S = S_i$, and $p_1(S_i) = \frac{n_i}{n_1}$

$$p(element2) = p_1(S_i) * (\frac{n_1}{n_2}) * \frac{1}{n_i} = \frac{1}{n_2}$$

If in $step_1$, the selected set $S = S_{d1}$, and $p_1(S_{d1}) = 1 - \frac{n_i}{n_1}$

$$p(element2) = (p_1(S_{d1}) + (1 - \frac{n_1}{n_2}) * p_1(S_i)) * \frac{1}{n_2 - n_i} = \frac{1}{n_2}$$

and

$$p(element1 = element2) = p_1(S_i) * (\frac{n_1}{n_2}) = \frac{n_i}{n_2}$$

# Case1: n1 > n2

$$(p_1(S_{d1}) = 1 - \frac{n_i}{n_1}) > (p_2 S_i = \frac{n_i}{n_2}) \tag{7}$$

So in $step_1$, when select $S_i$, it should be changed $S_{d2}$ in probability of $p$. The equation is

$$p_2(S_i) = p_1(S_i) * (1 - p) = \frac{n_i}{n_2}$$

$$p_2(S_{d2}) = p_2(S_{d1}) + p_1(S_i) * p = 1 - \frac{n_i}{n_2} \tag{8}$$

then

$$p = 1 - \frac{n_1}{n_2}$$

# Dependently Sampling Algorithm II

## Algorithm

- $step_1$: randomly select a set $S$ from $S_i$ and $S_{d1}$
- $S_1$:
    - $step_{12}$: if the $S$ is $S_{d1}$, randomly select a element1 from $S_{d1}$
    - $step_{13}$: if the $S$ is $S_i$, go to $step_3$
- $S_2$:
    - $step_{22}$: if the $S$ is $S_{d2}$, randomly select a element2 from $S_{d2}$
    - $step_{23}$: if the $S$ is $S_i$
        - randomly let $S = S_{d1}$ in probability of $\frac{n_1}{n_2}$
        - if $S$ is $S_{d2}$, randomly select a element2 from $S_{d2}$
        - if $S$ is $S_i$, go to $step_3$
- $step_3$: randomly select a element from $S_i$

# Dependently Sampling Algorithm II

## Probability

If in $step_1$, the selected set $S = S_i$, and $p_1(S_i) = \frac{n_i}{n_1}$

$$p(element2) = p_1(S_i) * (\frac{n_1}{n_2}) * \frac{1}{n_i} = \frac{1}{n_2}$$

If in $step_1$, the selected set $S = S_{d1}$, and $p_1(S_{d1}) = 1 - \frac{n_i}{n_1}$

$$p(element2) = (p_1(S_{d1}) + (1 - \frac{n_1}{n_2}) * p_1(S_i)) * \frac{1}{n_2 - n_i} = \frac{1}{n_2}$$

and

$$p(element1 = element2) = p_1(S_i) * (\frac{n_1}{n_2}) = \frac{n_i}{n_2}$$

# Solution 3
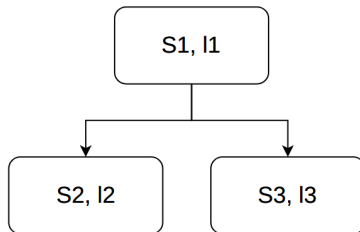
- First, We randomly select an idx $i_1$ from the $S_1$.
- If $n_1 < n_2$:
    - If $i_1 \notin S_i$, randomly sample in $S_2 - S_i$
    - If $i_1 \in S_i$ and $p > \frac{n_2}{n_1}$, randomly sample in $S_2 - S_i$
    - If $i_1 \in S_i$ and $p < \frac{n_2}{n_1}$, $i_2 = i_1$
- If $n_1 > n_2$:
    - If $i_1 \notin S_i$ and $p < \frac{n_1 * (n_2 - n_i)}{n_2 * (n_1 - n_i)}$, randomly sample in $S_2 - S_i$
    - If $i_1 \notin S_i$ and $p > \frac{n_1 * (n_2 - n_i)}{n_2 * (n_1 - n_i)}$, randomly sample in $S_i$
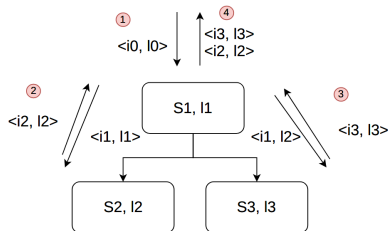    - If $i_1 \in S_i$, $i_2 = i_1$

# Sampling Tree

### Attributes

- There are two sets $S_a, S_b$
- $S_1 = S_a \cap S_b$
- $S_2 = S_a - S_b$
- $S_3 = S_b - S_a$
- $|S_2| < |S_3|$

# Sampling

- 1. if $p >= l_0/l_1$, sample $i_1$ from $S_1$. Otherwise $i_1 = i_0$
- 2. if $p >= l_1/l_2$, sample $i_2$ from $S_2$ and $i_1 = -1$. Otherwise $i_2 = i_1$
- 3. if if $i_1 \neq -1$ and $p >= l_2/l_3$, sample from $S_3$. Otherwise $i_3 = i_1$
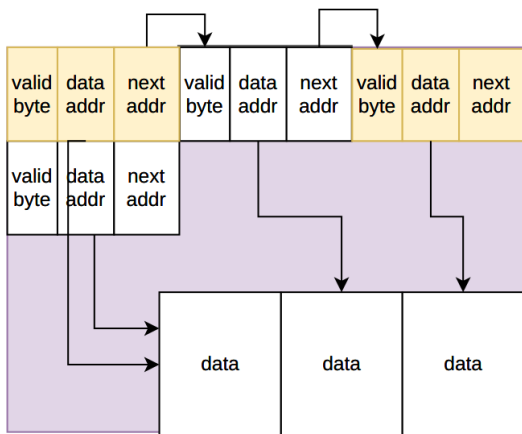- 4. return $< i_2, l_2 >, < i_3, l_3 >$

# Buffer Pool: Data Structure

### data

- There are two kinds of nodes: inode and datanode
- Every task has a head inode address
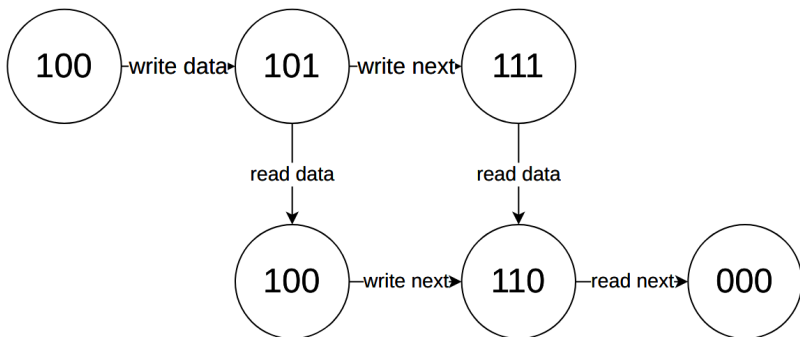
# Buffer Pool: Valid Byte

## valid byte

- data bit: If the data bit is equal to 1, the data addr is valid. Otherwise invalid

- next bit: If the next bit is equal to 1, the next addr is valid. Otherwise invalid

- used bit: If the used bit is equal to 1, this inode is used by some tasks

# Buffer Pool: Automata

### valid byte

| used bit | next bit | data bit |

# Buffer Pool: allocate inode

### case1

There is enough free space to allocate

```
1        if inode_tail + inode_size > data_head:
2            return inode_tail
```

### case2

Free Some unused inode

```
1        for head in all_heads:
2            if check_free(head) is True:
3                return head
```

# Buffer Pool: allocate data node

### case1

There is enough free space to allocate

```
1  if inode_tail + inode_size > data_head:
2      return inode_tail
```

### case2

Free Some unused datanode

```
1  free = True
2  for datanode in all_datanodes:
3      for ref in refs of datanode:
4          if databit(ref) == 0 && dataaddr(ref) == datanode:
5              free = False
6              break
7      if free is True:
8          return datanode
```
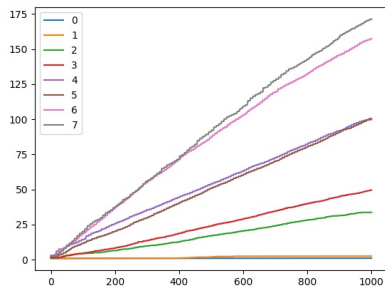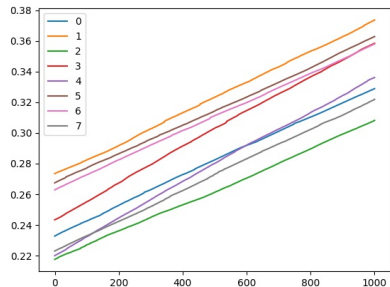
# Table of Contents

# Experiment



Figure: time



Figure: time with GlobalDataLoader