# FYP JC2104

## Presentation

**Imsong Jeon 1155087995**
**8 December 2021**

# 1. Objectives

FYP JC2104: Open Topics on Distributed Systems for Data Analytics

1. Self-Learning
2. Defining Problem and Solution

# 2. Process

2021-2022 Term 1

Summer | Beginning | Middle | End

# 2. Process

**Summer**





**Beginning**

Ursa

    C++

    CMake

    Cluster and Job Configuration

# 2. Process

## Middle

Example Codes

      hello_world.cc

      pi.cc

      word_count.cc

      kmeans.cc

      ...

from close point of view:

- High-level APIs
- Source Code

from far point of view:

- Workflow
- Structure

# 2. Process

## Middle

from close point of view:

- High-level APIs
- Source Code

```
100         auto data = TextSourceDataset(input, tg, n_partitions)
101                    .FlatMap([features](const std::string& line) { return LibsvmParse(line, features); })
102                    .PartitionBy(
103                        [n_partitions](const DataObj& vec) {
104                            std::random_device rd;
105                            std::mt19937 mt(rd());
106                            std::uniform_int_distribution<int> dist(0, n_partitions);
107                            return dist(mt);
108                        },
109                        n_partitions);
```

```
17    #include "base/properties.h"
18    #include "common/closure.h"
19    #include "common/dataset/dataset_partition.h"
20    #include "common/dataset/source_dataset.h"
21    #include "common/dataset/table_dataset.h"
22    #include "common/job_driver.h"
23    #include "common/resource_predictor.h"
24    #include "common/task_graph.h"
25
```

```
20
21    #include "common/constants.h"
22    #include "common/dataset/dataset.h"
23    #include "common/dataset/dataset_partition.h"
24    #include "common/io/input/csv_line_inputformat.h"
25    #include "common/io/input/hdfs_file_splitter.h"
26    #include "common/io/input/hdfs_input_block_info.h"
27    #include "common/io/input/line_inputformat.h"
28    #include "common/io/input/nfs_file_splitter.h"
29    #include "common/io/input/nfs_input_block_info.h"
30    #include "common/source_data.h"
31
```
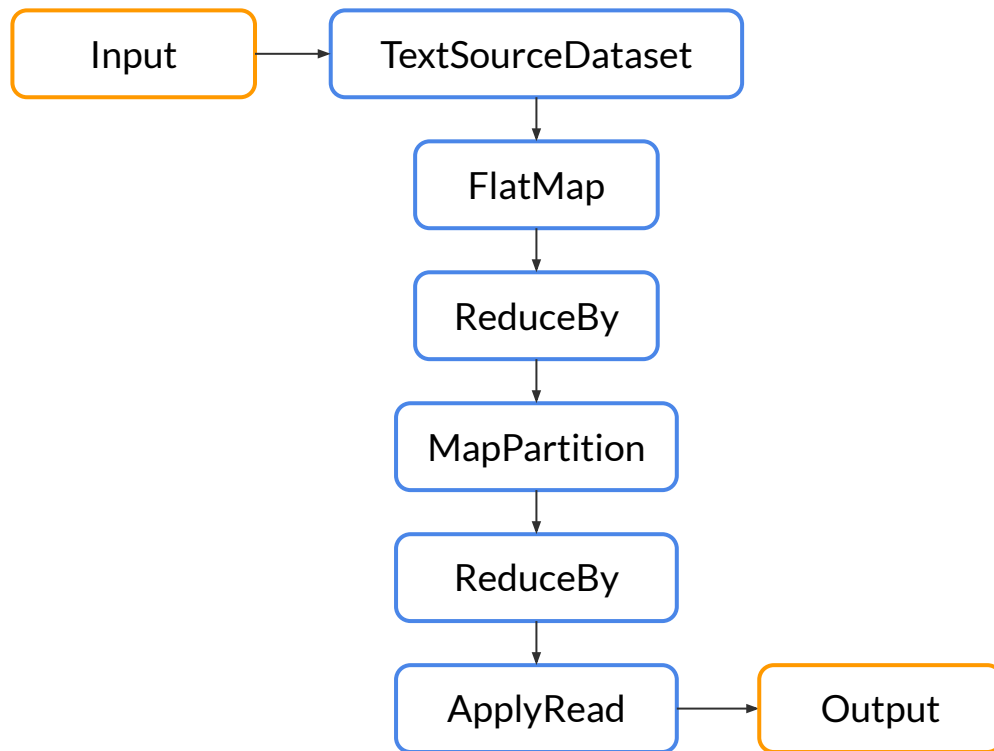
```
23
24    #include "common/engine.h"
25
```

kmeans.cc                 common/engine.h                              common/dataset/source_dataset.h

# 2. Process

## Middle

from far point of view:

- Workflow (DAG)
- Overall Structure

```
Input → TextSourceDataset
              ↓
           FlatMap
              ↓
           ReduceBy
              ↓
         MapPartition
              ↓
           ReduceBy
              ↓
         ApplyRead → Output
```

Workflow of word_count.cc

# 2. Process

**End**

More Example Codes

Implementations

    Word Length Count

    Parallel Dijkstra

# 3. Results

Word  Length Count

Input:
"I have two cats and two dogs"

Word Count output:

```
I1207 00:28:26.212761 182415 word_count.cc:62] 6
```

Word Length Count output:

```
I1207 00:29:46.871788 183663 word_length_count.cc:59] output: 1: 1
I1207 00:29:46.871800 183663 word_length_count.cc:59] output: 3: 3
I1207 00:29:46.871802 183663 word_length_count.cc:59] output: 4: 3
```

# 3. Results

Word Length Count/Run

```cpp
class WordCountJob : public Job {
 public:
  void Run(TaskGraph* tg, const std::shared_ptr<Properties>& config) const override {
    TextSourceDataset(config->Get("input"), tg, std::stoi(config->Get("parallelism")))
      //load data
        .FlatMap([](const std::string& line) {
          DatasetPartition<std::pair<int, int>> ret;
          ParseLine(ret, line);
          return ret; //dataset of (length, number of word length in this line)
        })
        .ReduceBy([](const std::pair<int, int>& ele) { return ele.first; },
                  [](std::pair<int, int>& agg, const std::pair<int, int>& update) { agg.second += update.second; }, 1)
        .ApplyRead([](auto data) {
          int count;
          count = data.size();
          for(int i = 0; i < count; i++){
            std::pair<int, int> cur_data = data.at(i);
            LOG(INFO) << "output: " << cur_data.first << ": " << cur_data.second;
          }
          google::FlushLogFiles(google::INFO);
        });
  }
};
```

# 3. Results

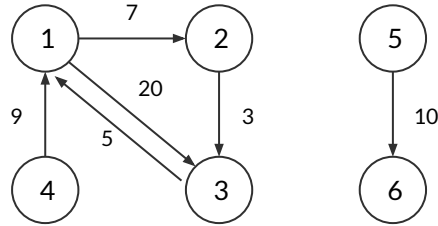Word Length Count/ParseLine

```cpp
27  void ParseLine(DatasetPartition<std::pair<int, int>>& collection, const std::string& line) {
28    if (line.empty()) {
29      return;
30    }
31    axe::base::WhiteSpaceTokenizer tokenizer(line);
32    std::string tok;
33    std::unordered_map<int, int> length_count;
34    while (tokenizer.next(tok)) {
35      length_count[tok.length()] += 1;
36    }
37    for (auto& pair : length_count) {
38      collection.push_back(pair);
39    }
40  }
41
```

# 3. Results

## Parallel Dijkstra

Input:

```
1    1 2 7
2    1 3 20
3    2 3 3
4    3 1 5
5    4 1 9
6    5 6 10
```



Configuration:
n-iterations: 2
start node: 1

Output:

```
I1207 00:48:46.596313 122074 parallel_dij_imsong.cc:225] NodeID = 1 Distance =  0 Parent NodeID =  -1
I1207 00:48:46.596325 122074 parallel_dij_imsong.cc:225] NodeID = 3 Distance =  10 Parent NodeID =  2
I1207 00:48:46.596328 122074 parallel_dij_imsong.cc:225] NodeID = 2 Distance =  7 Parent NodeID =  1
```

# 3. Results

Parallel Dijkstra/Node class

```cpp
28  class Node{
29    public:
30      Node() : neighbor_list_(std::make_shared<std::vector<std::pair<int,int>>>()) {}
31      Node(int node_id, int min_distance, int parent_node_id) :
32        node_id_(node_id), min_distance_(min_distance), parent_node_id_(parent_node_id),
33        neighbor_list_(std::make_shared<std::vector<std::pair<int,int>>>()) {}
34      Node(int node_id, int min_distance, int parent_node_id, const std::shared_ptr<std::vector<std::pair<int,int>>>& neighbor_list) :
35        node_id_(node_id), min_distance_(min_distance), parent_node_id_(parent_node_id), neighbor_list_(neighbor_list) {}
36
37      int GetNodeID() const { return node_id_; }
38      int GetMinDistance() const { return min_distance_; }
39      int GetParentNodeID() const { return parent_node_id_; }
40      const std::shared_ptr<std::vector<std::pair<int,int>>>& GetNeighborList() const { return neighbor_list_; }
41
42      double GetMemory() const {
43        double ret = sizeof(int) * 3 + neighbor_list_->size() * sizeof(std::pair<int,int>);
44        return ret;
45      }
46      bool operator<(const Node& other) const { return node_id_ < other.node_id_; }
47      bool operator==(const Node& other) const { return node_id_ == other.node_id_; }
48
49      friend void operator<<(axe::base::BinStream& bin_stream, const Node& n) {
50        bin_stream << n.node_id_ << n.min_distance_ << n.parent_node_id_ << *(n.neighbor_list_); }
51      friend void operator>>(axe::base::BinStream& bin_stream, Node& n) {
52        bin_stream >> n.node_id_ >> n.min_distance_ >> n.parent_node_id_ >> *(n.neighbor_list_); }
53
54    private:
55      int node_id_;
56      int min_distance_;
57      int parent_node_id_;
58      std::shared_ptr<std::vector<std::pair<int,int>>> neighbor_list_;
59  };
```

# 3. Results

Parallel Dijkstra/Run (part 1)

```cpp
110   class ParallelDijkstra : public Job {
111    public:
112     void Run(TaskGraph* tg, const std::shared_ptr<Properties>& config) const override {
113       auto input = config->GetOrSet("data", "");
114       int n_partitions = std::stoi(config->GetOrSet("parallelism", "20"));
115       int n_iters = std::stoi(config->GetOrSet("n_iters", "3"));
116       int start_ID = std::stoi(config->GetOrSet("start_node", "1"));
117       // Load data
118       auto graph = TextSourceDataset(input, tg, n_partitions)
119                         .FlatMap([start_ID](const std::string& line) { return ParseLine(line, start_ID); })
120                         .ReduceBy([](const Node& ele) { return ele.GetNodeID(); },
121                           [](Node& agg, const Node& ele) {
122                             std::shared_ptr<std::vector<std::pair<int,int>>> aggList = agg.GetNeighborList();
123                             std::shared_ptr<std::vector<std::pair<int,int>>> eleList = ele.GetNeighborList();
124                             aggList->push_back(eleList->front());
125                           },
126                           n_partitions);
127
```

```
I1207 11:46:07.515542 256070 parallel_dij_imsong.cc:104] Node ID = 1 minDist = 0 parent_node_id = -1 neighbors = 2 7 3 20
```

# 3. Results

Parallel Dijkstra/Run (part 2)

```cpp
172    auto updated_graph = std::make_shared<axe::common::Dataset<Node>>(graph);
173    std::shared_ptr<axe::common::Dataset<std::pair<int, std::pair<int, int>>>> min_distances;
174
175    for(int i = 0; i < n_iters; ++i){
176      min_distances = std::make_shared<axe::common::Dataset<std::pair<int, std::pair<int, int>>>>(updated_graph->MapPartition(get_distances)
177                           .ReduceBy([](const std::pair<int, std::pair<int, int>>& ele) { return ele.first; },
178                             [](std::pair<int, std::pair<int, int>>& agg, const std::pair<int, std::pair<int, int>>& ele) {
179                               if(agg.second.first > ele.second.first){
180                                 agg.second.first = ele.second.first; //update min distance
181                                 agg.second.second = ele.second.second; //update parent ID
182                               }
183                             },
184                             n_partitions));
185
186      updated_graph = std::make_shared<axe::common::Dataset<Node>>(
187          updated_graph->SharedDataMapPartitionWith(min_distances.get(), apply_updates));
188    }
```

```cpp
139    auto get_distances = [](const DatasetPartition<Node>& graph){
140      DatasetPartition<std::pair<int, std::pair<int, int>>> distances;
141      for(const Node& n : graph){
142        distances.push_back(std::make_pair(n.GetNodeID(), std::make_pair(n.GetMinDistance(), n.GetParentNodeID())));
143        for(std::pair<int,int> neighbor : *n.GetNeighborList()){
144            int new_distance = MAX_INT;
145            if(n.GetMinDistance() < MAX_INT)
146              new_distance = neighbor.second + n.GetMinDistance();
147            distances.push_back(std::make_pair(neighbor.first, std::make_pair(new_distance, n.GetNodeID())));
148        }
149      }
150      return distances;
151    };
```

# 3. Results

Parallel Dijkstra/Run (part 3)

```cpp
153    auto apply_updates = [](const DatasetPartition<Node>& graph, const DatasetPartition<std::pair<int, std::pair<int, int>>>& min_distances) {
154        DatasetPartition<Node> updated_graph;
155        for(const std::pair<int, std::pair<int, int>>& p : min_distances){
156            int node_id = p.first;
157            bool updated = false;
158            for(const Node& n : graph){
159                if(n.GetNodeID() == node_id){
160                    updated_graph.push_back(Node(node_id, p.second.first, p.second.second, n.GetNeighborList()));
161                    updated = true;
162                    break;
163                }
164            }
165            if(!updated){
166                updated_graph.push_back(Node(node_id, p.second.first, p.second.second));
167            }
168        }
169        return updated_graph;
170    };
171
```

```cpp
197    auto get_output = [](const DatasetPartition<Node>& graph){
198        DatasetPartition<std::pair<int, std::pair<int, int>>> output;
199        for(const Node& n : graph){
200            if(n.GetMinDistance() < MAX_INT){
201                output.push_back(std::make_pair(n.GetNodeID(), std::make_pair(n.GetMinDistance(), n.GetParentNodeID())));
202            }
203        }
204        return output;
205    };
206
207    updated_graph->MapPartition(get_output)
208                 .PartitionBy([](const std::pair<int, std::pair<int, int>>&) { return 0; }, 1)
209                 .ApplyRead( [](const auto& partition) {
210                         for (auto& par : partition) {
211                             LOG(INFO) << "NodeID = " << par.first << " Distance =  " << par.second.first << " Parent NodeID =  " << par.second.second;
212                         }
```

# 4. Plans and Ideas

Plans:

- Implement and test with larger data and more workers
- More study
- Distributed PCA (Principal Component Analysis)

Ideas

- Linear Algebra Package
- Real-time interaction

# 5. Conclusion

Thank you!