

#### Fast CSV Loading Using GPUs and RDMA for In-Memory Data Processing

BTW 2021

Alexander Kumaigorodski, Clemens Lutz, Volker Markl *TU Berlin, DFKI* 



NVLink, RDMA, GPUDirect, ... IO with 100+ Gbit/s



Efficiently load and parse CSV files end-to-end







Parsing on GPUs

- Parsing on GPUs
- Streamed loading with GPUDirect RDMA

- Parsing on GPUs
- Streamed loading with GPUDirect RDMA
- Evaluate fast interconnects for end-to-end streamed loading

# Parsing on GPUs

## **Background: Warp Divergence**

- No branch prediction on GPUs
- Threads execute in groups of 32, a *warp*
- Warp execute same instruction during cycle
- if(cond) {...} else {...} can cut performance in half

#### **Problem: CSV Parsing**

## **Problem: CSV Parsing**

ID, Country, ISOCode, Population\n 1, Germany, DE, 83149319\n 2, Spain, ES, 47007367\n 3, France, FR, 67076431\n 4, Italy, IT, 60317116\n 5, Japan, JP, 126150745\n 6, China, CN, 1427647786\n 7, United States of America, US, 328239523\n 8, Canada, CA, 37894799\n

#### Parsing on GPUs

Parsers typically have complex control flow (warp divergence)

#### Parsing on GPUs

- Parsers typically have complex control flow (warp divergence)
- Exploring new trade-off: Simplify control flow at expense of additional data passes (700+ GB/s)













#### **Conceptual Overview: Fast Mode**



#### **Conceptual Overview: Quoted Mode**



### Main Challenges

- Partitioning the data into chunks for parallel processing
- Determining each chunk's context
- Deserializing fields in parallel



#### Load Balancing Warps









ChunkDelimiterschunk 0ID, Country, ISOCode, Population 1, Germany, DE, 83149319 2, Spain, ES, 411chunk17007367 n3, France, FR, 67076431 14, Italy, IT, 60317116 n5, Japan, JP, 126112chunk 250745 n6, China, CN, 1427647786 n7, United States of America, US, 3282398chunk 366

523\n8, Canada, CA, 37894799\n9, Australia, AU, 25721892\n10, Russian F...

•••



•••



Chunk Delimiters Prefix sum (offset) chunk 0 ID, Country, ISOCode, Population 1, Germany, DE, 83149319 2, Spain, ES, 4 0 11 chunk1 7007367vn3, France, FR, 67076431vn4, Italy, IT, 60317116vn5, Japan, JP, 1261 12 11 chunk 2 50745<u>6</u>, China, CN, 1427647786<u>7</u>, United States of America, US, 328239 8 23 chunk 3 523\n8, Canada, CA, 37894799\n9, Australia, AU, 25721892\n10, Russian F... 10 31

•••

•••

•••





#### **Fast Mode: Quoted Field Delimiters**

ID,Name,Philosophy\n
1,"Aristotle","Quality is not an act, it is a habit."\n
2,"Plato","When men speak ill of thee, live so as nobody may believe them."\n
3,"Epictetus","It's not what happens to you, but how you react to it that matters."\n
...

#### **Fast Mode: Quoted Field Delimiters**

ID,Name,Philosophy\n
1,"Aristotle","Quality is not an act, it is a habit."\n
2,"Plato","When men speak ill of thee, live so as nobody may believe them."\n
3,"Epictetus","It's not what happens to you, but how you react to it that matters."\n
...

ID	Name	Philosophy
1	"Aristotle"	"Quality is not an act, it is a habit."
2	"Plato"	"When men speak ill of thee, live so as nobody may believe them."
3	"Epictetus"	"It's not what happens to you, but how you react to it that matters."
		•••

#### **Fast Mode: Quoted Field Delimiters**

ID,Name,Philosophy\n
1,"Aristotle","Quality is not an act, it is a habit."\n
2,"Plato","When men speak ill of thee, live so as nobody may believe them."\n
3,"Epictetus","It's not what happens to you, but how you react to it that matters."\n
...

ID	Name	Philosophy
1	"Aristotle"	"Quality is not an act, it is a habit."
2	"Plato"	"When men speak ill of thee, live so as nobody may believe them."
3	"Epictetus"	"It's not what happens to you, but how you react to it that matters."
•••	•••	

ID	Name	Philosophy
1	"Aristotle"	"Quality is not an act
it is a habit."	2	"Plato"
"When men speak ill of thee	live so as nobody may believe them."	3
	••••	••••
### **Quoted Mode**

Character is considered *quoted* whenever the number of preceding quotation marks is uneven.

### **Quoted Mode**

ID,Name,Philosophy\n
1,"Aristotle","Quality is not an act, it is a habit."\n
2,"Plato","When men speak ill of thee, live so as nobody may believe them."\n
3,"Epictetus","It's not what happens to you, but how you react to it that matters."\n
...

#### FieldsIndex

Θ	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	stream	0	1	2	3	4	5	6	7	8	9	10	11
0	3	8	19	21	33	0	73	75	0	149	151	0	163	0	233		0	3	8	19	21	33	73	75	149	151	163 2	.33



Deserialization challenges on GPU:

- Parallelization
- Warp divergence and occupancy



- Deserialization challenges on GPU:
  - Parallelization
  - Warp divergence and occupancy
- Row-based: complex kernels, inefficient memory access



- Deserialization challenges on GPU:
  - Parallelization
  - Warp divergence and occupancy
- Row-based: complex kernels, inefficient memory access
- Column-based: same data type, same kernel logic and control flow



- Like SQL's DDL, users specify a column's max length along its type char(15) int(8) float(12) ...
- Every thread deserializes one field, warp 32 fields in parallel









# Streaming

## Streaming







#### **Context Handover**













# **Evaluation**

### **Experiment Setup**

- Intel Xeon Gold 5115 (10c/20t)
- Nvidia Tesla V100 (PCIe 3.0)
- Nvidia Tesla V100 (NVLink 2.0)
- InfiniBand 100 Gb/s

 NYC Yellow Taxi Trips (1.9 GB, 22.5M rows, 18 cols)

TPC-H Lineitem (719 MB, 6M rows, 16 cols)

• int\_444
 (1 GB, 70M rows, 3 columns)

## NYC Yellow Taxi 30 GB/s 25 GB/s 20 GB/s INdH9N02HT 15 GB/s 10 GB/s 5 GB/s 0 GB/s

#### NYC Yellow Taxi

	30 GB/s									
	25 GB/s									
HPUT	20 GB/s									
THROUGH	15 GB/s	max, bandwidth	PCIe 3.0							
	10 GB/s									
	5 GB/s						2.52 GB/s			
	0.05 (	0.03 GB/s	0.08 GB/s	0.17 GB/s		0.25 GB/s				
	U UB/S	PostgreSQL (CPU)	HyPer (CPU)	OmniSci (CPU)	-	csvmonkey (CPU)	InstantLoading (CPU-32c)			



#### NYC Yellow Taxi





















#### Fast Mode vs. Quoted Mode



#### Fast Mode vs. Quoted Mode



### **Hardware Scalability**



### **Hardware Scalability**


## **Hardware Scalability**



## **Hardware Scalability**



## **Hardware Scalability**



## Conclusion

• GPUs can solve data loading bottleneck

• New way to integrate GPUs into DBs