






chunked reading and memory consumption in columnflow


Mathis Frahm

columnflow meeting

29.11.2024

chunked reading in Creat... 														
Task	Events	Chunksize	#	Pool size	debug?	materialization	 Runtime (total)	 Runtime (loop body)	#	Max memory (VIRT)	#	Max memory (RES)		
CreateHistogramr	2271259	1k		1	yes	PARTITIONS				6800		4700		
		1k		1	no	PARTITIONS				5600		3500		
		1k		2	no	PARTITIONS				6000		3800		
		1k		20	no	PARTITIONS				15500		9300		
		1k		1	yes	SLICES				3607		1611		
		1k		1	no	SLICES				3991		1615		
		1k		20	no	SLICES				5245		1620		
		50k		2	no	PARTITIONS	01:58	01:10		4067		1615		
		100k		2	no	PARTITIONS	01:12	00:50		4083		1638		
		200k		2	no	PARTITIONS	00:59	00:46		4100		1676		
		200k		4	no	PARTITIONS	01:00	00:43		4240		1668		
		50k		2	no	SLICES (fixed)	01:13	00:54		4100		1651		
		100k		2	no	SLICES (fixed)	00:54	00:43		4100		1662		
		200k		2	no	SLICES (fixed)	00:51	00:44		4100		1707		
		200k		4	no	SLICES (fixed)	00:49	00:42		4400		1715		
		400k		4	no	SLICES (fixed)	00:45	00:42		4400		1810		

- Input files: reduced events (500mb) + producer outputs (1400mb)
- Chunk size and pool size (after fix) does not significantly affect memory consumption (we always load full parquet files)
- larger chunks and more pools improve runtime
- processing during histogramming (complicated expressions) might introduce memory spikes during processing

chunked reading in Select... 															
Task	Events	Chunksize	#	Pool size	debug?	materialization	⌚ Runtime (total)	⌚ Runtime (loop body)	#	Max memory (VIRT)	#	Max memory (RES)			
SelectEvents	1042230	50k	2		no	SLICES (fixed)	02:32	01:26		3000		1400			
		100k	1		no	SLICES (fixed)	02:25	01:09		3200		1650			
		100k	2		no	SLICES (fixed)	01:56	01:06		4100		1850			
		100k	4		no	SLICES (fixed)	01:56	01:06		4600		2300			
		200k	2		no	SLICES (fixed)	01:50	00:59		3800		2300			
		50k	2		no	PARTITIONS	02:40	01:29		3100		1400			
		100k	2		no	PARTITIONS	02:19	01:08		3400		1850			
		100k	4		no	PARTITIONS	02:03	01:08		4200		2200			
		200k	2		no	PARTITIONS	01:46	01:00		4300		1950			
		200k	4		no	PARTITIONS	01:42	01:02		4700		2950			
		400k	1		no	PARTITIONS	02:34	00:55		5200		3300			
		400k	2		no	PARTITIONS	01:36	00:55		5600		3500			
		> filesize	1		no	PARTITIONS	02:22	01:03		8100		5600			
		200k	2		yes	PARTITIONS	02:29	00:56		3200		1900			
		100k	2		yes	PARTITIONS	02:47	01:02		2900		1500			

- pool size of 2 seems to be optimal (no runtime improvements with >2, significant runtime increase with ==1)
- chunksize of 100k as compromise between runtime and memory consumption (<2GB for 2.2GB nano input file)
 - reducing to 50k increases runtime by 30% and reduces memory consumption by 25%
 - increasing to 200k improves runtime by 5% and increases memory consumption by 25%
 - might be very different for other analyses/selectors (e.g. memory spikes during processing of chunk)