



Energy-Efficient Sorting using Solid State Disks

The Sort Benchmark	Algorithms	Solid State Disks																																																							
<p>The Benchmark</p> <ul style="list-style-type: none"> Sort 100 byte records with a 10 byte key Introduced 1985; starting with ~100 MB New categories added targeting Speed/Size/Throughput (GraySort) Time (MinuteSort) Cost Efficiency (PennySort) Energy Efficiency (JouleSort, 2007) <ul style="list-style-type: none"> 10 GB, 100 GB, 1000 GB 100 TB (2010) Classes: Indy (tuned), Daytona (general) 	<p>External Memory Multiway Mergesort</p> <ul style="list-style-type: none"> Phase 1: Run Formation Phase 2: Merge Runs Careful parameter selection for optimal performance while requiring a single merge pass Parallel implementations utilize the 4 CPU threads Overlapping of I/O and computation Run Formation uses key extraction and radixsort Two implementations: <p>EcoSort (Indy: 10 GB, 100 GB)</p> <ul style="list-style-type: none"> Bring overlapping to the limits Allow independent tuning of more parameters <p>DEMsort (Indy: 1000 GB, 100 TB)</p> <ul style="list-style-type: none"> Developed by Sanders, Singler et al. at the Karlsruhe Institute of Technology Won the 2009 Sort Benchmark in the categories MinuteSort and GraySort using a 200-node cluster Efficient also on a single node Allows in-place sorting, needed to sort 1000 GB with just 1024 GB of storage <p>Nsort (Daytona: 100 GB, 1000 GB)</p> <ul style="list-style-type: none"> Commercial software Sorts arbitrary data types 	<p>Pro</p> <ul style="list-style-type: none"> Built from NAND flash memory chips No mechanically moving parts Good shock resistance Low energy consumption Higher throughput than HDDs <p>Con</p> <ul style="list-style-type: none"> Higher price and less capacity than today's HDDs Small block random writes are slow Performance may degrade depending on access pattern Properties vary depending on manufacturer, model, firmware: 																																																							
<p>2007</p> <p>Rivoire, Shah, Ranganathan, Kozyrakis Stanford University and HP Labs</p>	<p>2010</p> <p>Beckmann, Meyer, Sanders, Singler Goethe University and Karlsruhe Institute of Technology</p>																																																								
<p>2010</p> <p>Beckmann, Meyer, Sanders, Singler Goethe University and Karlsruhe Institute of Technology</p>	<p>2011 (to be submitted)</p>																																																								
<p>2007</p> <p>Intel Core 2 Duo T7600 (Mobile CPU) 2 cores, 2 threads, 1.66 GHz</p> <p>2 PCI-e Disk Controllers (8+4 SATA) 1 SATA (onboard) 13 x Hitachi Travelstar 5K160 160 GB Notebook HDD</p> <p>XFS on Linux Software Raid (Striping) NSort (commercial sorter)</p> <p>Power Idle 59 W Power Loaded 100 W</p> <p>2007 JouleSort Winner 10 GB, 100 GB</p>	<p>2010</p> <p>Intel Atom 330 2 cores, 4 threads, 1.6 GHz</p> <p>4 GB 4 x SATA 3.0 Gb/s (onboard)</p> <p>4 x SuperTalent FTM56GX25H 256 GB SSD</p> <p>Linux XFS on Linux Software Raid (Striping) EcoSort, DEMsort using STXXL</p> <p>25 W 37 W</p> <p>2010 Sort Benchmark, first 100 TB results to be submitted 2011.</p>	<p>Results</p> <p>Winner of the 2010 Sort Benchmark in the JouleSort categories Indy 10 GB, 100 GB and 1000 GB and Daytona 100 GB!</p> <table border="1"> <thead> <tr> <th rowspan="2">Class, Size [GB]</th> <th colspan="3">2007</th> <th colspan="3">2010</th> </tr> <tr> <th>Time [s]</th> <th>Energy [kJ]</th> <th>Rec./J</th> <th>Time [s]</th> <th>Energy [kJ]</th> <th>Rec./J</th> </tr> </thead> <tbody> <tr> <td>Indy, 10</td> <td>86.6</td> <td>8.6</td> <td>11628</td> <td>72.4</td> <td>2.3</td> <td>42635</td> </tr> <tr> <td>Indy, 100</td> <td>881</td> <td>88.1</td> <td>11354</td> <td>691</td> <td>25.1</td> <td>39853</td> </tr> <tr> <td>Daytona, 100</td> <td>881</td> <td>88.1</td> <td>11354</td> <td>756</td> <td>27.9</td> <td>35789</td> </tr> <tr> <td>Indy, 1000</td> <td>7196*</td> <td>2920*</td> <td>3425</td> <td>17026</td> <td>572</td> <td>17489</td> </tr> <tr> <td>Daytona, 1000</td> <td>7196*</td> <td>2920*</td> <td>3425</td> <td>1897*</td> <td>5273</td> <td>1.5</td> </tr> <tr> <td>Indy, 100 TB</td> <td>-</td> <td>-</td> <td>-</td> <td>9835**</td> <td>694 MJ**</td> <td>1441</td> </tr> </tbody> </table> <p>Using low power hardware does not imply an increase in running time: in the 10 GB and 100 GB category we beat previous results both in terms of energy consumption and running time. As a consequence of winning all three categories using a single machine, a new 100 TB JouleSort category was introduced for the 2010 Sort Benchmark, first 100 TB results to be submitted 2011.</p> <p>* regular server hardware, not a low energy machine ** 200-node cluster</p>	Class, Size [GB]	2007			2010			Time [s]	Energy [kJ]	Rec./J	Time [s]	Energy [kJ]	Rec./J	Indy, 10	86.6	8.6	11628	72.4	2.3	42635	Indy, 100	881	88.1	11354	691	25.1	39853	Daytona, 100	881	88.1	11354	756	27.9	35789	Indy, 1000	7196*	2920*	3425	17026	572	17489	Daytona, 1000	7196*	2920*	3425	1897*	5273	1.5	Indy, 100 TB	-	-	-	9835**	694 MJ**	1441
Class, Size [GB]	2007			2010																																																					
	Time [s]	Energy [kJ]	Rec./J	Time [s]	Energy [kJ]	Rec./J																																																			
Indy, 10	86.6	8.6	11628	72.4	2.3	42635																																																			
Indy, 100	881	88.1	11354	691	25.1	39853																																																			
Daytona, 100	881	88.1	11354	756	27.9	35789																																																			
Indy, 1000	7196*	2920*	3425	17026	572	17489																																																			
Daytona, 1000	7196*	2920*	3425	1897*	5273	1.5																																																			
Indy, 100 TB	-	-	-	9835**	694 MJ**	1441																																																			
<p>2011</p> <p>Intel Atom 330 2 cores, 4 threads, 1.6 GHz</p> <p>4 GB 4 x SATA 3.0 Gb/s (onboard)</p> <p>4 x SuperTalent FTM56GX25H 256 GB SSD</p> <p>Linux XFS on Linux Software Raid (Striping) EcoSort, DEMsort using STXXL</p> <p>25 W 37 W</p> <p>2010 Sort Benchmark, first 100 TB results to be submitted 2011.</p>	<p>I/O and CPU utilization while sorting 10 GB:</p>	<p>Results</p> <p>Winner of the 2010 Sort Benchmark in the JouleSort categories Indy 10 GB, 100 GB and 1000 GB and Daytona 100 GB!</p> <table border="1"> <thead> <tr> <th rowspan="2">Class, Size [GB]</th> <th colspan="3">2007</th> <th colspan="3">2010</th> </tr> <tr> <th>Time [s]</th> <th>Energy [kJ]</th> <th>Rec./J</th> <th>Time [s]</th> <th>Energy [kJ]</th> <th>Rec./J</th> </tr> </thead> <tbody> <tr> <td>Indy, 10</td> <td>86.6</td> <td>8.6</td> <td>11628</td> <td>72.4</td> <td>2.3</td> <td>42635</td> </tr> <tr> <td>Indy, 100</td> <td>881</td> <td>88.1</td> <td>11354</td> <td>691</td> <td>25.1</td> <td>39853</td> </tr> <tr> <td>Daytona, 100</td> <td>881</td> <td>88.1</td> <td>11354</td> <td>756</td> <td>27.9</td> <td>35789</td> </tr> <tr> <td>Indy, 1000</td> <td>7196*</td> <td>2920*</td> <td>3425</td> <td>17026</td> <td>572</td> <td>17489</td> </tr> <tr> <td>Daytona, 1000</td> <td>7196*</td> <td>2920*</td> <td>3425</td> <td>1897*</td> <td>5273</td> <td>1.5</td> </tr> <tr> <td>Indy, 100 TB</td> <td>-</td> <td>-</td> <td>-</td> <td>9835**</td> <td>694 MJ**</td> <td>1441</td> </tr> </tbody> </table> <p>Using low power hardware does not imply an increase in running time: in the 10 GB and 100 GB category we beat previous results both in terms of energy consumption and running time. As a consequence of winning all three categories using a single machine, a new 100 TB JouleSort category was introduced for the 2010 Sort Benchmark, first 100 TB results to be submitted 2011.</p> <p>* regular server hardware, not a low energy machine ** 200-node cluster</p>	Class, Size [GB]	2007			2010			Time [s]	Energy [kJ]	Rec./J	Time [s]	Energy [kJ]	Rec./J	Indy, 10	86.6	8.6	11628	72.4	2.3	42635	Indy, 100	881	88.1	11354	691	25.1	39853	Daytona, 100	881	88.1	11354	756	27.9	35789	Indy, 1000	7196*	2920*	3425	17026	572	17489	Daytona, 1000	7196*	2920*	3425	1897*	5273	1.5	Indy, 100 TB	-	-	-	9835**	694 MJ**	1441
Class, Size [GB]	2007			2010																																																					
	Time [s]	Energy [kJ]	Rec./J	Time [s]	Energy [kJ]	Rec./J																																																			
Indy, 10	86.6	8.6	11628	72.4	2.3	42635																																																			
Indy, 100	881	88.1	11354	691	25.1	39853																																																			
Daytona, 100	881	88.1	11354	756	27.9	35789																																																			
Indy, 1000	7196*	2920*	3425	17026	572	17489																																																			
Daytona, 1000	7196*	2920*	3425	1897*	5273	1.5																																																			
Indy, 100 TB	-	-	-	9835**	694 MJ**	1441																																																			