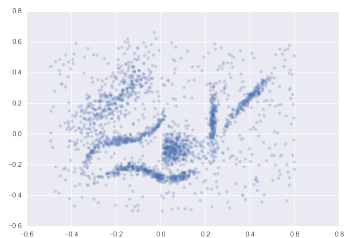


Theoretically Efficient Parallel Density-Peaks Clustering

Michael Huang
Under the direction of
Shangdi Yu and Prof. Julian Shun

October 16, 2022

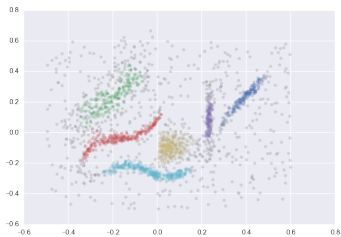
Density-based Clustering



Unclustered data



k-means clustering result¹

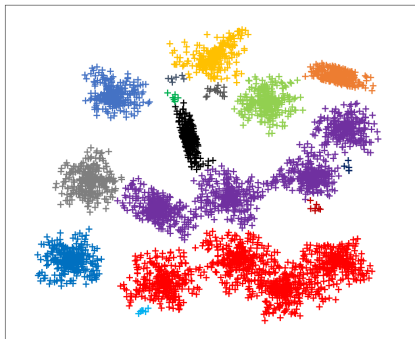


DBSCAN clustering result²

¹Everitt, Landau, and Leese 2009.

²Ester et al. 1996.

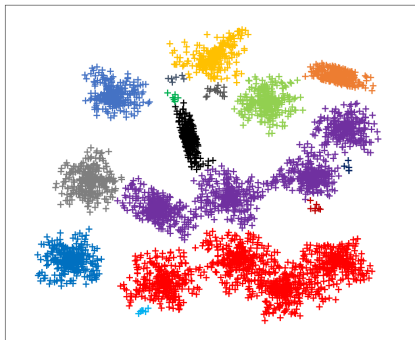
Why Density-Peaks Clustering (DPC) Algorithm



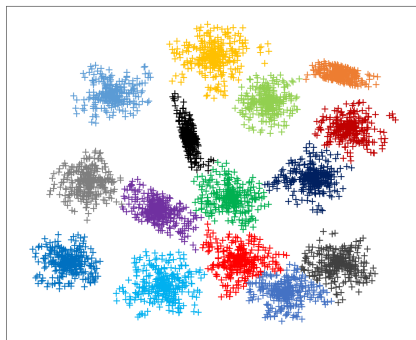
DBSCAN fails on datasets where clusters are close together³

³Amagata and Hara 2021.

Why Density-Peaks Clustering (DPC) Algorithm



DBSCAN fails on datasets where clusters are close together³

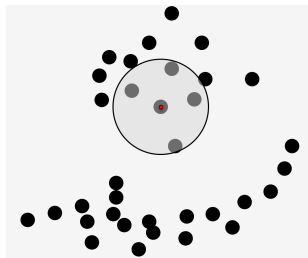


DPC is able to separate close clusters⁴

³Amagata and Hara 2021.

⁴Amagata and Hara 2021.

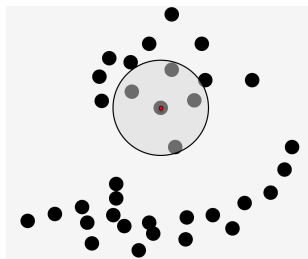
DPC Algorithm Procedure Description



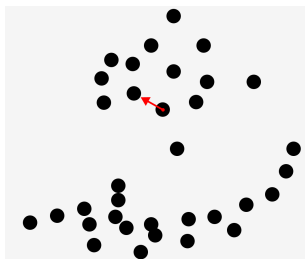
Compute density⁵

⁵Rodriguez and Laio 2014.

DPC Algorithm Procedure Description



Compute density⁵

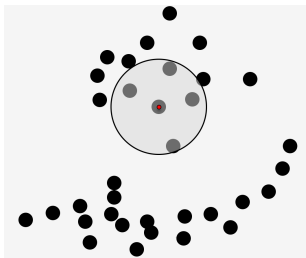


Find dependent point
(the nearest neighbor
with higher density)⁶

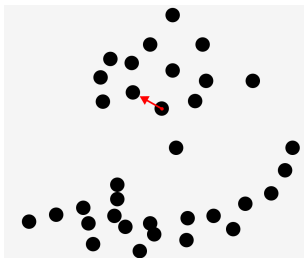
⁵Rodriguez and Laio 2014.

⁶Rodriguez and Laio 2014.

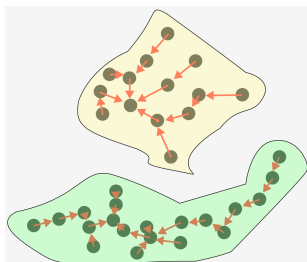
DPC Algorithm Procedure Description



Compute density⁵



Find dependent point
(the nearest neighbor
with higher density)⁶



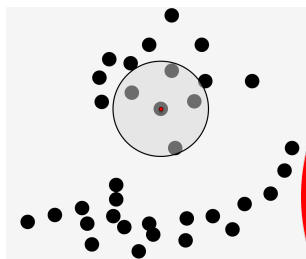
Separate into clusters⁷

⁵Rodriguez and Laio 2014.

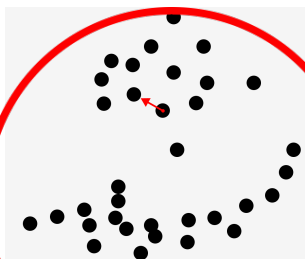
⁶Rodriguez and Laio 2014.

⁷Rodriguez and Laio 2014.

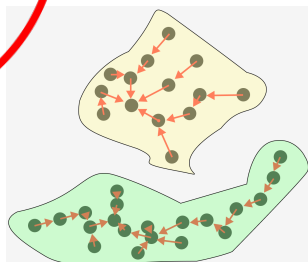
DPC Algorithm Procedure Description



Compute density⁵



Find dependent point
(the nearest neighbor
with higher density)⁶



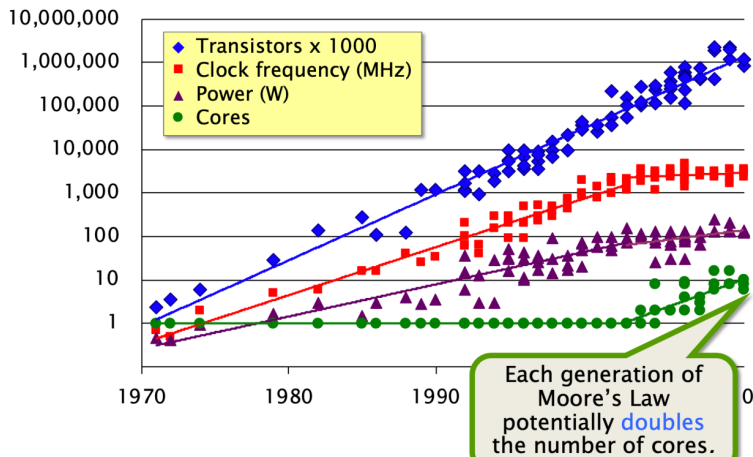
Separate into clusters⁷

⁵Rodriguez and Laio 2014.

⁶Rodriguez and Laio 2014.

⁷Rodriguez and Laio 2014.

Why Focus on Parallelism



© 2019 Julian Shun Slide adapted from 6.172 (Charles Leiserson and Suman Amarasinghe)

CPU clock-speed hits ceiling; #cores increases exponentially⁸

⁸Shun 2021.

Parallel Algorithm Background

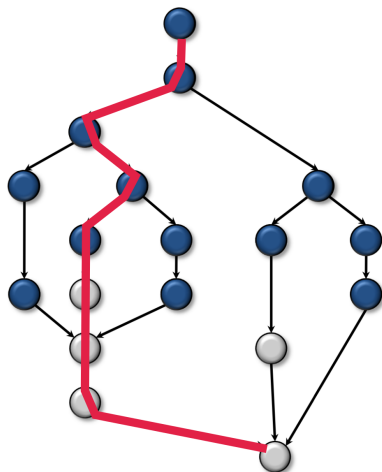
T_p = runtime with p processors

T_1 = work

T_∞ = span

Brent's Law:

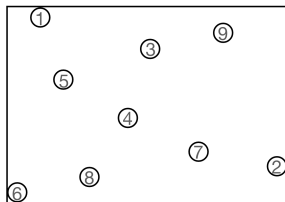
$$T_p \leq T_\infty + \frac{T_1 - T_\infty}{p}$$



Computational graph of a parallel algorithm⁹

⁹Shun 2021.

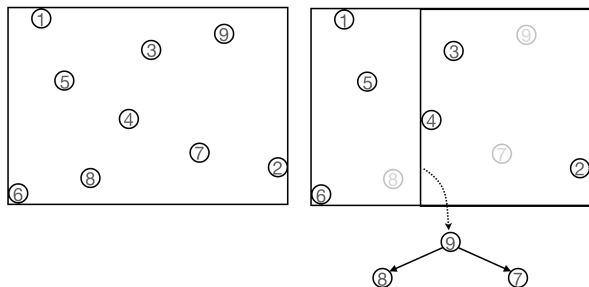
Parallel Dependent Point Finding with Priority Search kd-tree



Binary Space Partitioning Tree:

- 1 Divide points up equally
- 2 Satisfy heap property (higher in the tree \implies higher density)

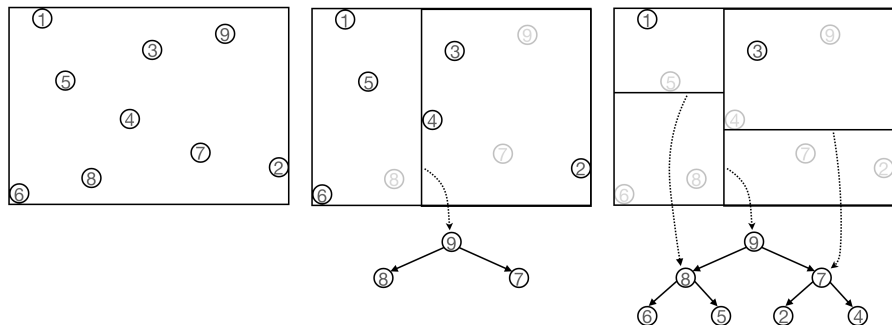
Parallel Dependent Point Finding with Priority Search kd-tree



Binary Space Partitioning Tree:

- 1 Divide points up equally
- 2 Satisfy heap property (higher in the tree \implies higher density)

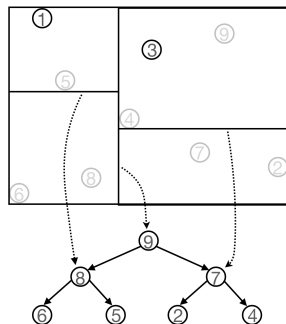
Parallel Dependent Point Finding with Priority Search kd-tree



Binary Space Partitioning Tree:

- 1 Divide points up equally
- 2 Satisfy heap property (higher in the tree \implies higher density)

Parallel Dependent Point Finding with Priority Search kd-tree

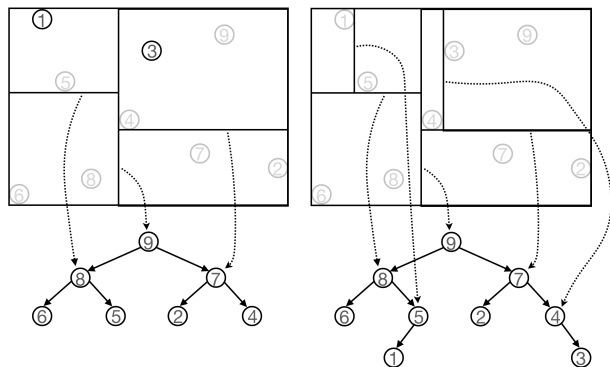


Binary Space Partitioning Tree:

- 1 Divide points up equally
- 2 Satisfy heap property (higher in the tree \implies higher density)

Reduce dependent point finding
from $O(n)$ to Avg. $O(\log(n))$

Parallel Dependent Point Finding with Priority Search kd-tree

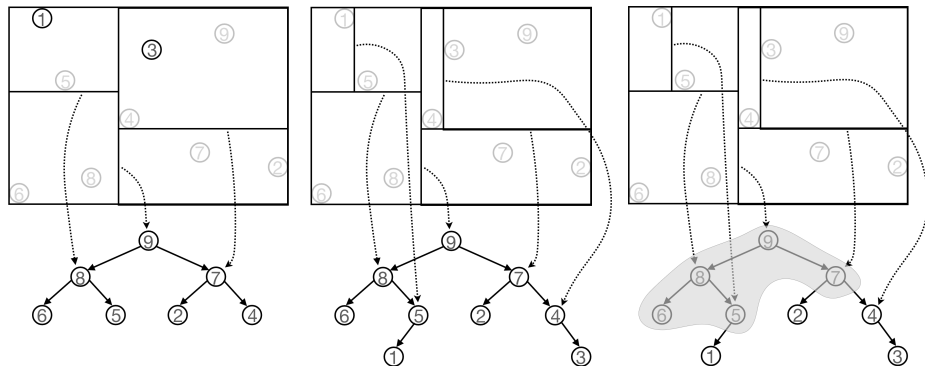


Binary Space Partitioning Tree:

- 1 Divide points up equally
- 2 Satisfy heap property (higher in the tree \implies higher density)

Reduce dependent point finding
from $O(n)$ to Avg. $O(\log(n))$

Parallel Dependent Point Finding with Priority Search kd-tree



Binary Space Partitioning Tree:

- 1 Divide points up equally
- 2 Satisfy heap property (higher in the tree \implies higher density)

Reduce dependent point finding
from $O(n)$ to Avg. $O(\log(n))$

Algorithmic Complexity

Algorithms	Compute density		Find dependent point	
	Avg. Work	Avg. Span	Avg. Work	Avg. Span
Previous SOTA ¹⁰	$O(n^{2-\frac{1}{d}} + n\rho)$	$O(n^{1-\frac{1}{d}} + \rho)$	$O(n^2)$	$O(n)$
Our algorithm	$O(n^{2-\frac{1}{d}})$	$O(n^{1-\frac{1}{d}})$	$O(n \log(n))$	$O(\log(n))$

Complexity comparison

- 1 n : the number of points to be clustered
- 2 ρ : average density of points
- 3 d : the number of dimensions each point has

¹⁰Rodriguez and Laio 2014; Amagata and Hara 2021.

Experiment Setup

- 1 30-core, 2-way hyperthreading, CPU @3.1 GHz
- 2 Implemented with ParlayLib¹¹ and ParGeo¹²



Google Cloud Platform

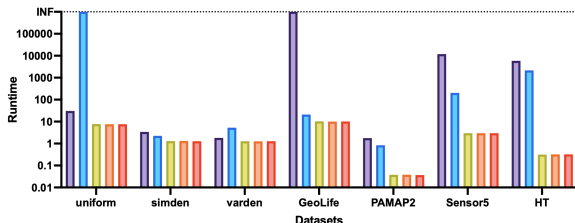
Dataset	n	d	synthetic
<i>uniform</i>	10M	2	yes
<i>simden</i>	10M	2	yes
<i>varde</i>	10M	2	yes
<i>GeoLife</i>	24.88M	3	no
<i>PAMAP2</i>	0.26M	4	no
<i>Sensor</i>	3.84M	5	no
<i>HT</i>	0.93M	8	no

Datasets

¹¹Blelloch, Anderson, and Dhulipala 2020.

¹²Wang et al. 2022.

Runtime Comparison

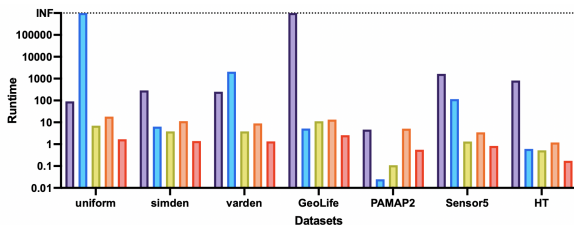


Overall speedups:
8.3–4666.3x

Density computation

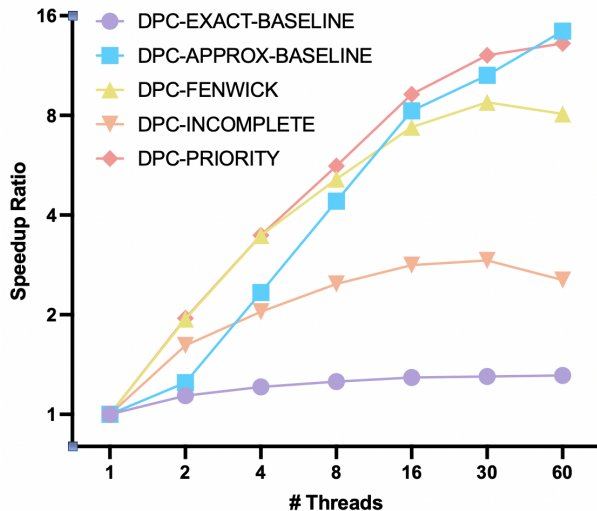
Algorithms

- DPC-EXACT-BASELINE
- DPC-APPROX-BASELINE
- DPC-FENWICK
- DPC-INCOMPLETE
- DPC-PRIORITY



Dependent point finding

Parallel Scalability



13.2x self-relative speedup

Conclusion

- 1 Proposed the Priority Search *kd*-tree data structure and proved its avg. query complexity
- 2 Developed a theoretically efficient and practically fast DPC algorithm, with up to 4666x speedup compared to SOTA

Acknowledgements

- Shangdi Yu
- Prof. Julian Shun
- MIT PRIMES Program