

Ultrametric Cluster Hierarchies: I Want 'em All!

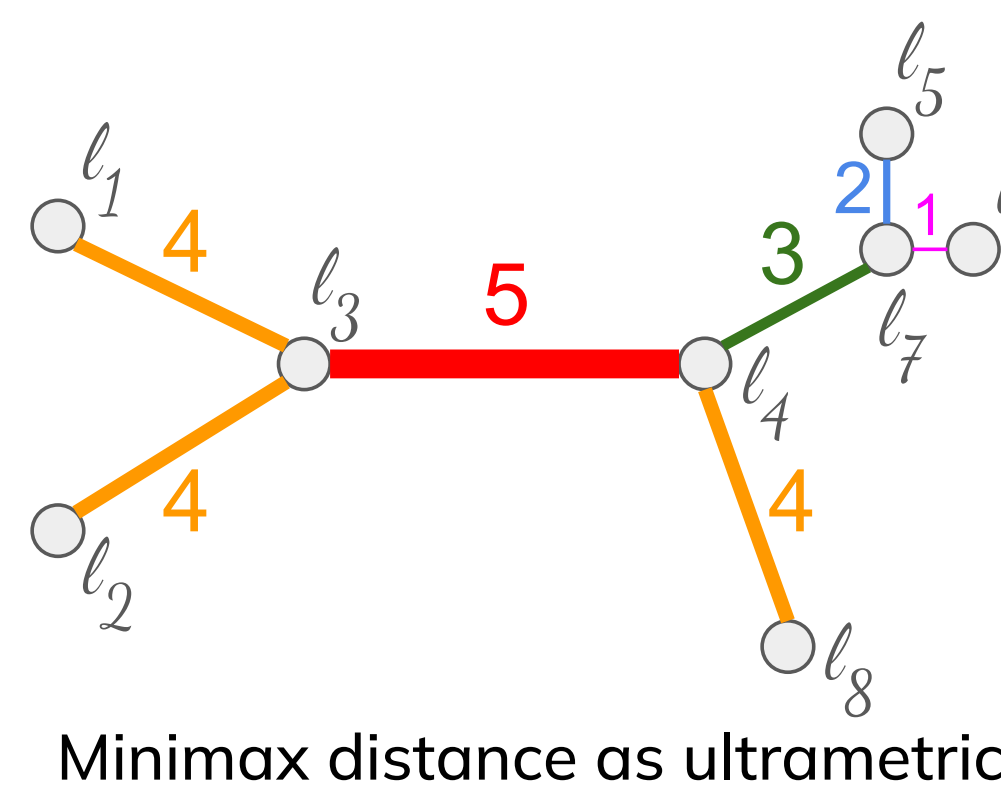
Andrew Draganov^{* 1,2}, Pascal Weber^{* 3,4,5},
Rasmus Jørgensen¹, Anna Beer³, Claudia Plant^{3,5}, Ira Assent^{1,2}

Similarity

All ultrametrics are hierarchical

An ultrametric is a metric which also satisfies the *strong triangle inequality*: $d(x, z) \leq \max \{d(x, y), d(y, z)\}$

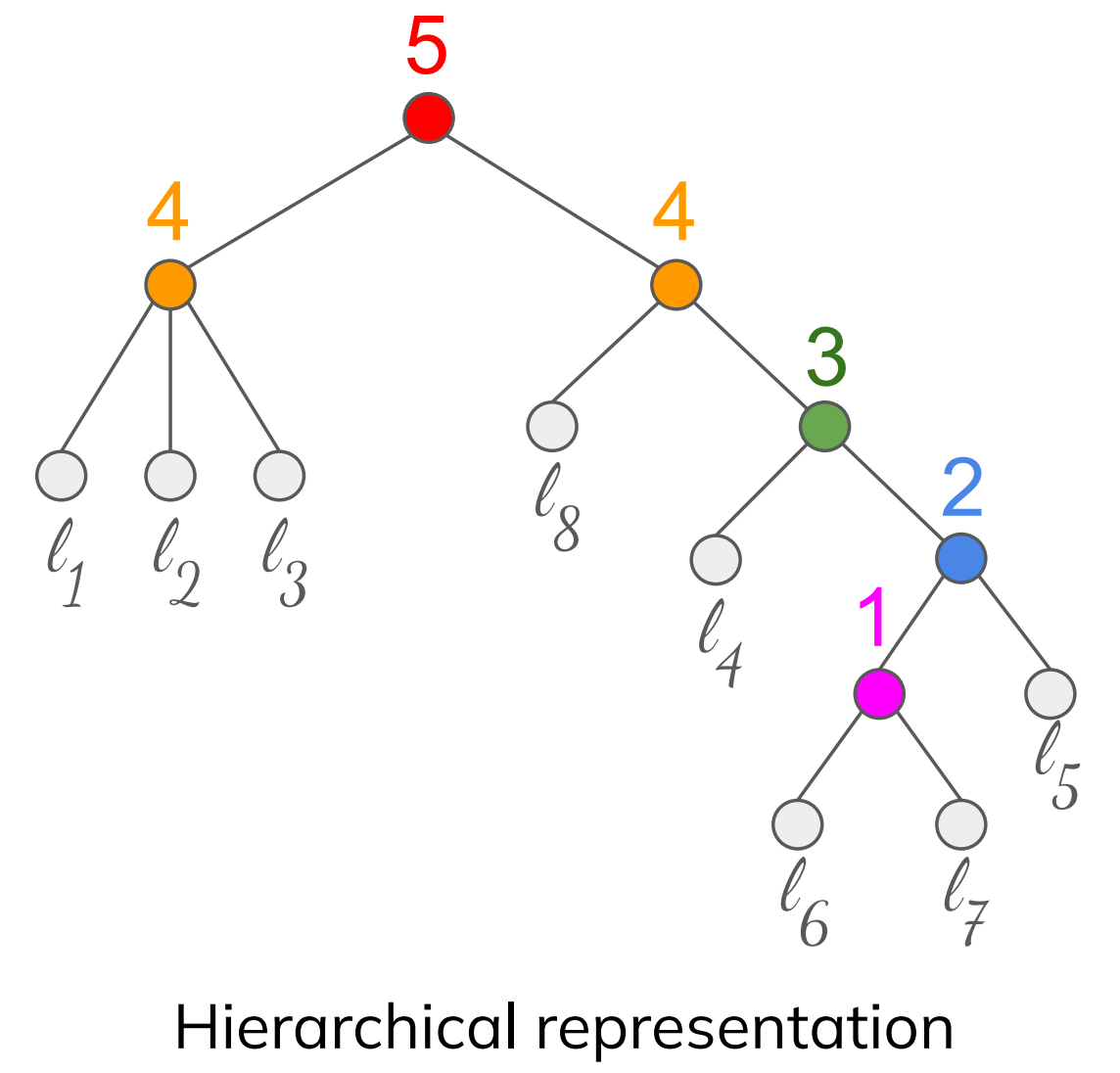
Every ultrametric can be represented by a **hierarchy**.



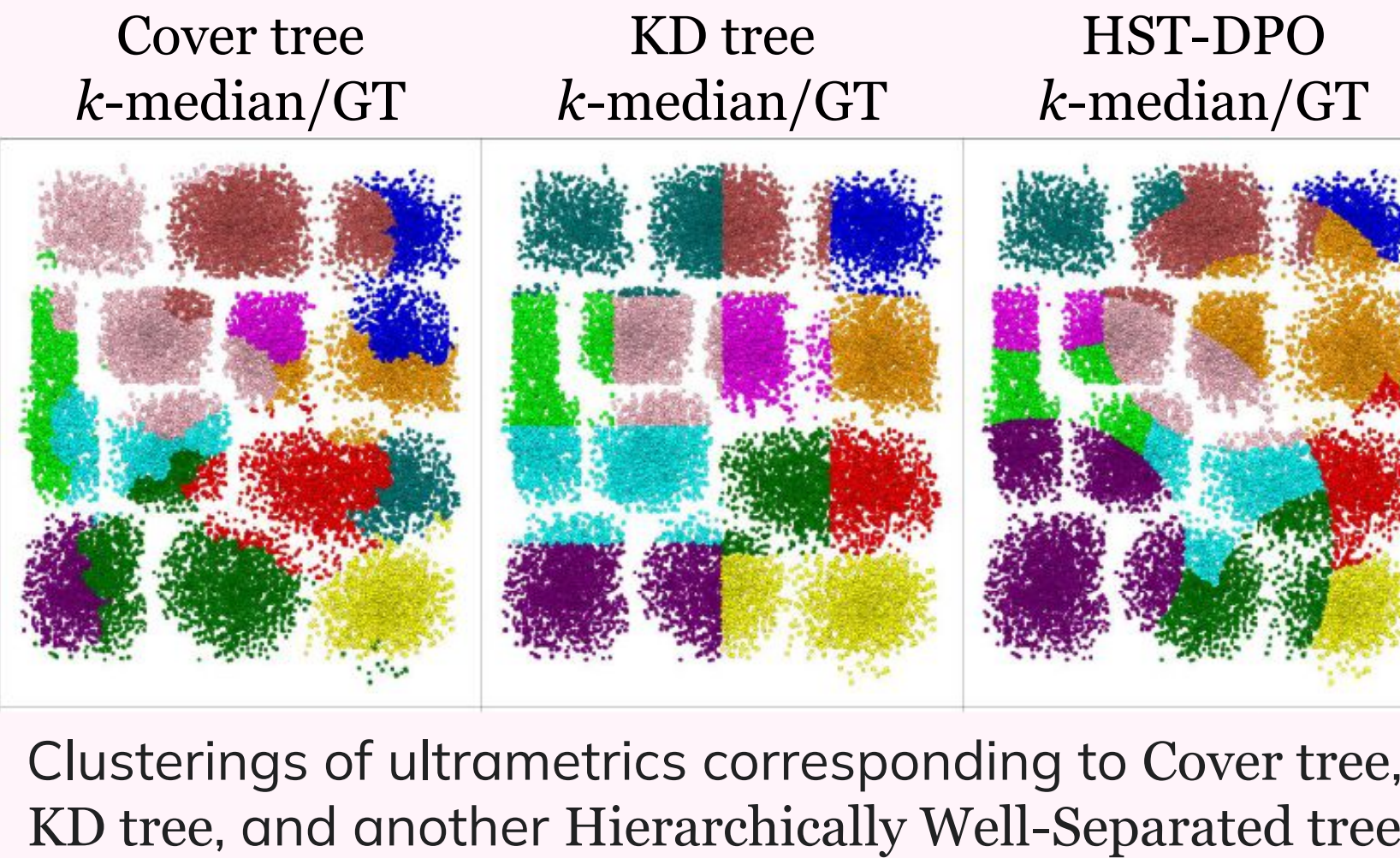
Distance matrix of minimax ultrametric

	l_1	l_2	l_3	l_4	l_5	l_6	l_7	l_8
l_1	0	4	4	5	5	5	5	5
l_2	4	0	4	5	5	5	5	5
l_3	4	4	0	5	5	5	5	5
l_4	5	5	5	0	3	3	3	4
l_5	5	5	5	3	0	2	2	4
l_6	5	5	5	3	2	0	1	4
l_7	5	5	5	3	2	1	0	4
l_8	5	5	5	4	4	4	4	0

Every hierarchy with *node values growing along the paths from the leaves to the root* corresponds to an ultrametric.



The distance between two nodes is the node value of their lowest common ancestor.



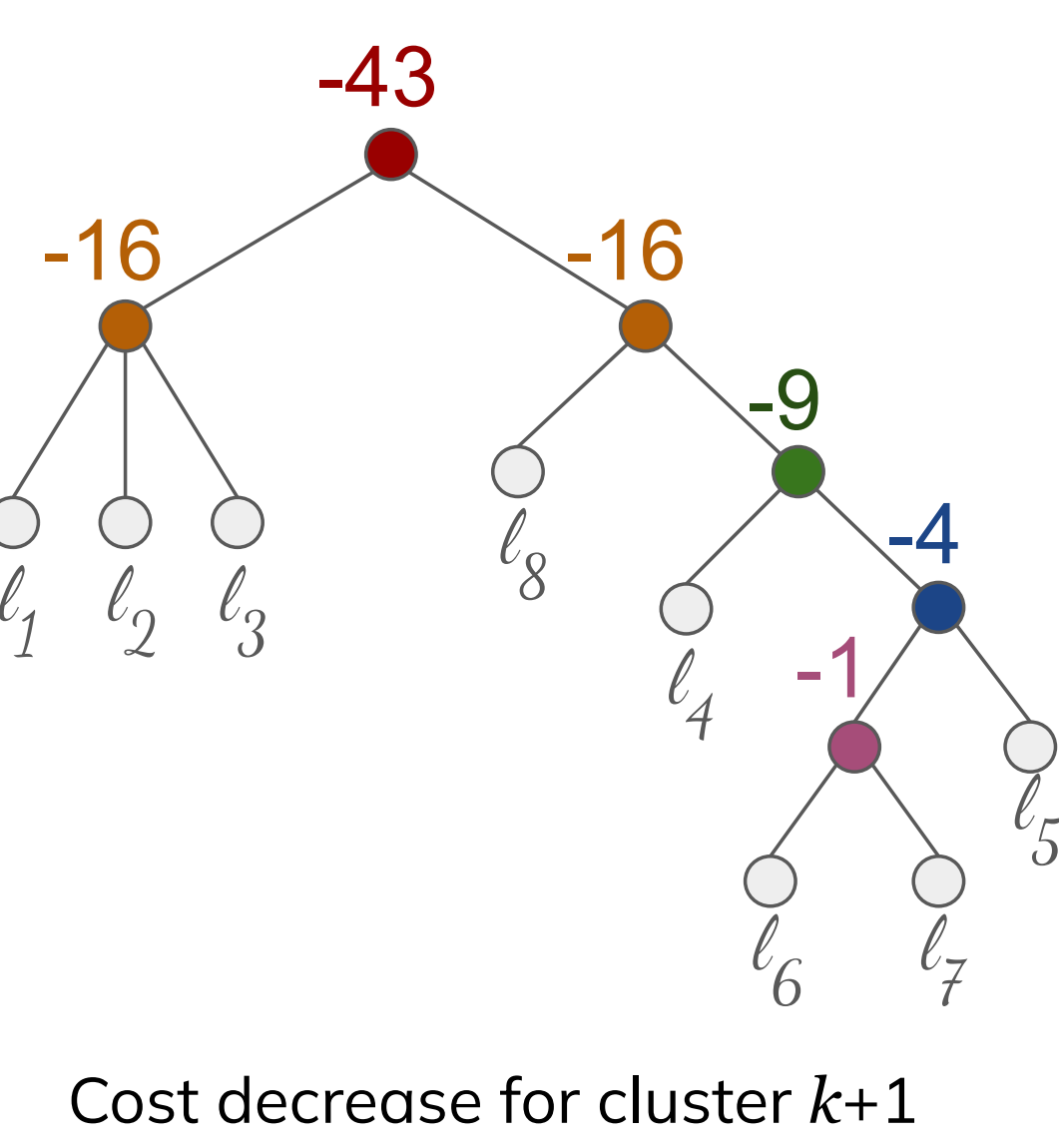
Cover trees, KD trees, and other Hierarchically Well-Separated trees can all represent data, but they all focus on different characteristics.

Hierarchy

Centroid-based clustering in ultrametrics

It takes **Sort(n)** time to find the optimal k -z (median, mean, etc.) solutions for **all** values of k in an ultrametric.

The **optimal solutions** are themselves hierarchical.



Cost decrease for cluster $k+1$

Clusterings	Total costs
$k = 1: l_1 l_2 l_3 l_4 l_5 l_6 l_7 l_8$	$105 = (3 \times 5^2 + 4^2 + 3^2 + 2^2 + 1^2)$
$k = 2: l_1 l_2 l_3 l_4 l_5 l_6 l_7 l_8$	$62 = (2 \times 4^2) + (4^2 + 3^2 + 2^2 + 1^2)$
$k = 3: l_1 l_2 l_3 l_4 l_5 l_6 l_7 l_8$	$46 = (4^2) + (4^2 + 3^2 + 2^2 + 1^2)$
$k = 4: l_1 l_2 l_3 l_4 l_5 l_6 l_7 l_8$	$30 = (4^2 + 3^2 + 2^2 + 1^2)$
$k = 5: l_1 l_2 l_3 l_4 l_5 l_6 l_7 l_8$	$14 = (3^2 + 2^2 + 1^2)$
$k = 6: l_1 l_2 l_3 l_4 l_5 l_6 l_7 l_8$	$5 = (2^2 + 1^2)$
$k = 7: l_1 l_2 l_3 l_4 l_5 l_6 l_7 l_8$	$1 = (1^2)$
$k = 8: l_1 l_2 l_3 l_4 l_5 l_6 l_7 l_8$	0

Cluster centers in bold | The sum of k -means costs

Partitioning

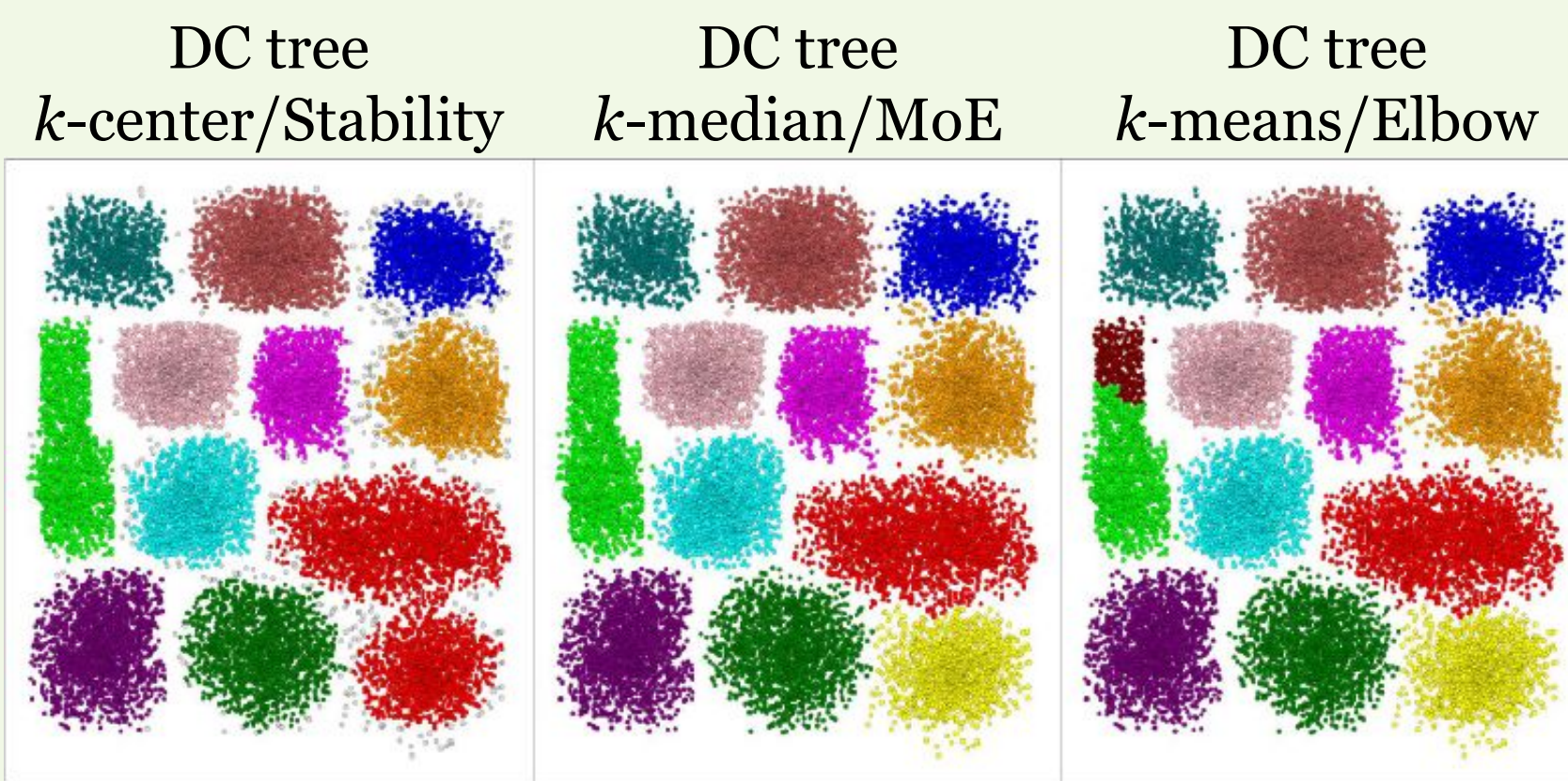
Extract a clustering from a hierarchy

Partitioning a hierarchy can be achieved extremely fast (**O(n)** time). Possible strategies are:

- Threshold the values in the tree (DBSCAN)
- Pick the “best” clustering by a function (HDBSCAN)
- Optimal clustering for a user-specified value of k
- Elbow method

DC tree captures density connectivity.

- Centroid **hierarchies** can split density-connected clusters
- Some **partitioning** methods find noise or determine the number of clusters automatically



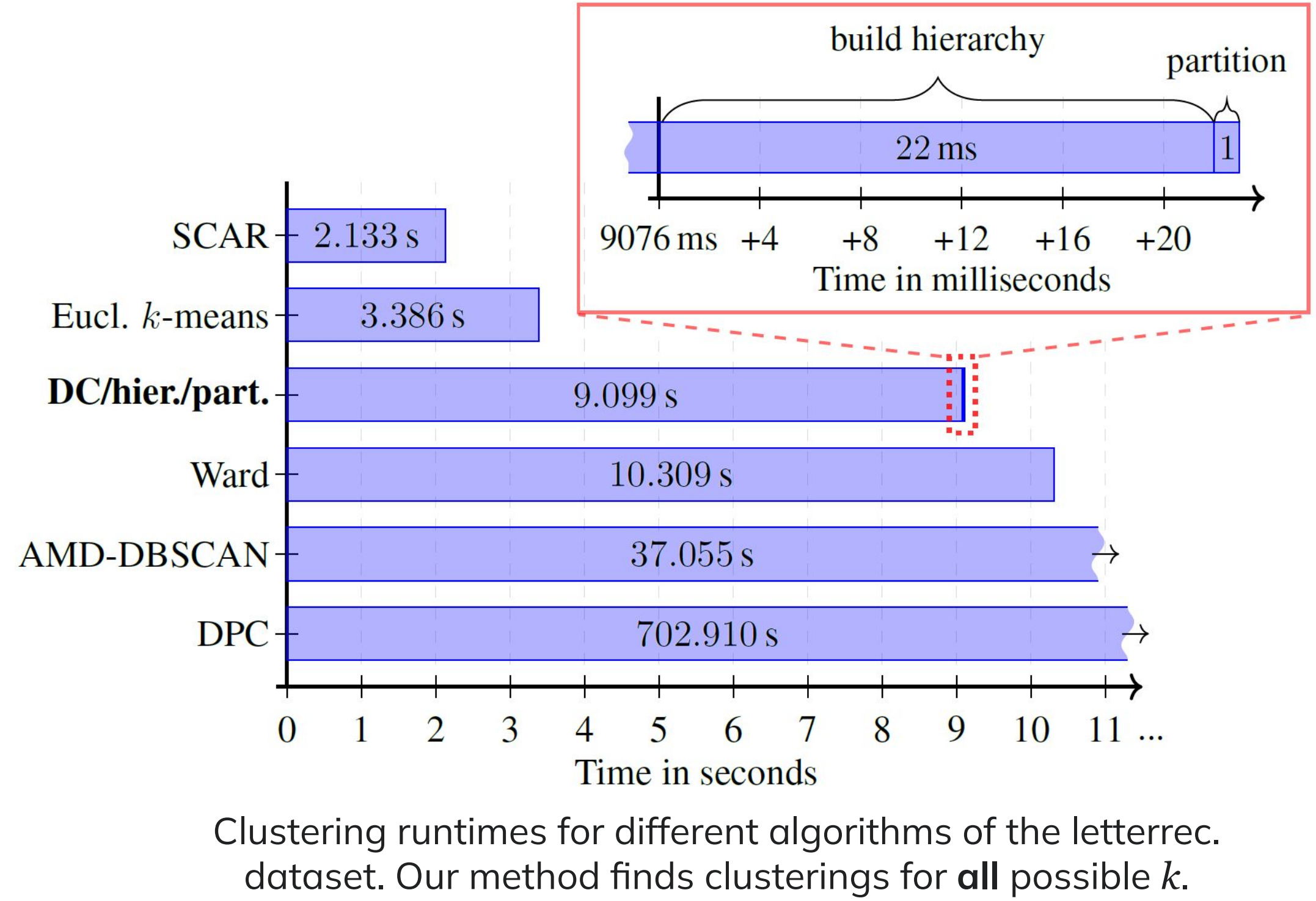
Clusterings when using the DC tree as similarity, and different centroid hierarchies and partitioning methods

Runtimes

High efficiency through single ultrametric preprocessing

Our framework requires only a **single upfront ultrametric computation**, after that we can generate **multiple different clusterings** with **negligible additional runtime**.

Faster methods compute the clustering for only one single parameter setting.



Clustering Quality

Flexibility enables deeper insights

There is no combination of hierarchy and partitioning that generally works “best”. Although DC tree/ k -means/Elbow performs good in many cases, it is not universally superior to the other combinations; different pairings excel on different datasets. Our framework allows to rapidly switch between different hierarchies and partitioning methods.

Dataset	k -center Stability	DC tree			k -center Stability	Cover tree			Eucl. k -means	competitors			
		k -median MoE	k -means Elbow	Elbow		k -median MoE	k -means Elbow	Elbow		SCAR	Ward	AMD-DBSCAN	DPC
Tabular Data	Boxes	90.1	99.3	97.9	2.6	42.1 ± 4.7	24.2 ± 1.6	24.2 ± 1.6	93.5 ± 4.3	0.1 ± 0.1	95.8	63.9	25.9
	D31	79.7	42.7	82.9	46.5 ± 1.8	62.0 ± 5.4	67.7 ± 3.2	67.7 ± 3.2	92.0 ± 2.7	41.7 ± 5.4	92.0	86.4	18.5
	airway	38.0	65.9	58.8	0.8	18.2 ± 2.4	12.0 ± 1.4	12.0 ± 1.4	39.9 ± 2.0	-0.9 ± 0.5	43.7	31.7	65.1
	lactate	41.0	41.0	67.5	0.1	4.1 ± 0.6	1.7 ± 0.2	1.7 ± 0.2	28.6 ± 1.1	1.5 ± 1.0	27.7	71.5	0.0
	HAR	30.0	46.9	52.8	14.7 ± 8.8	14.2 ± 4.7	9.6 ± 2.2	9.6 ± 2.2	46.0 ± 4.5	5.5 ± 3.2	49.1	0.0	33.2
	letterrec	12.1	16.6	17.9	5.8 ± 0.2	7.2 ± 0.6	6.2 ± 0.3	6.2 ± 0.3	12.9 ± 0.6	0.4 ± 0.1	14.7 ± 0.9	7.9	0.0
Image Data	PenDigits	66.4	73.1	75.4	8.0 ± 0.8	12.0 ± 0.6	8.9 ± 0.5	8.9 ± 0.5	55.3 ± 3.2	0.9 ± 0.3	55.2	55.6	28.8 ± 1.1
	COIL20	81.2	72.8	72.6	46.4 ± 4.4	46.6 ± 2.1	47.7 ± 2.0	47.7 ± 2.0	58.2 ± 2.8	33.5 ± 2.0	68.6	39.2	35.9 ± 0.1
	COIL100	80.1	66.8	70.0	44.6 ± 4.2	46.6 ± 1.5	50.1 ± 1.2	50.1 ± 1.2	56.1 ± 1.4	16.7 ± 0.8	61.4	14.2	0.2
	cmu_faces	60.2	56.6	66.5	8.6 ± 3.1	37.1 ± 4.1	34.2 ± 2.1	34.2 ± 2.1	53.2 ± 4.7	38.5 ± 2.9	61.6	0.7	0.6
	OptDigits	55.3	77.0	77.0	40.9 ± 3.5	20.9 ± 2.3	18.1 ± 2.4	18.1 ± 2.4	61.3 ± 6.6	14.4 ± 4.1	74.6 ± 2.4	63.2	0.0
	USPS	33.7	29.3	29.3	12.0 ± 1.7	8.7 ± 1.0	11.2 ± 1.5	11.2 ± 1.5	52.3 ± 1.7	2.9 ± 0.9	63.9	0.0	21.0
	MNIST	19.7	41.7	46.0	11.1 ± 1.7	5.4 ± 0.6	5.4 ± 0.6	5.4 ± 0.6	36.9 ± 1.0	1.3 ± 0.4	52.7	0.0	-

ARI values for the SHiP framework on the DC tree and Cover tree ultrametrics and competitors. Euclidean k -means, SCAR, and Ward are given the ground-truth value k .

References

Anna Beer, Andrew Draganov, et al. (2023). “Connecting the dots – density-connectivity distance unifies DBSCAN, k -center and Spectral Clustering”. In: SIGKDD Conference on Knowledge Discovery and Data Mining: p. 80–92.

Yuxiang Zeng, Yongxin Tong, and Lei Chen. (2021). “HST+: An efficient index for embedding arbitrary metric spaces”. In: IEEE 37th International Conference on Data Engineering: p. 648–659.



Our proposed **SHiP** clustering framework
(1) **fits an ultrametric (similarity)**
(2) **computes a centroid-based hierarchy**
(3) **extracts a partitioning**

- ✓ C++ Code on Github
- ✓ pip package of the Python interface

