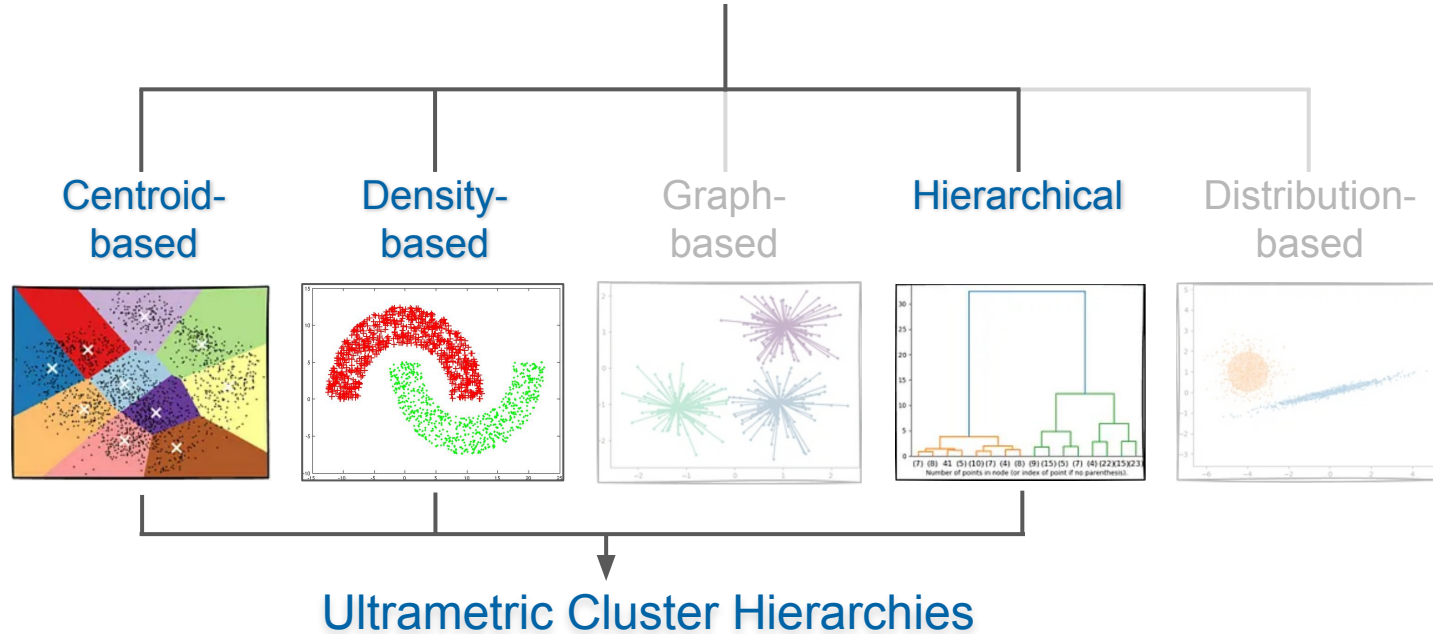# Ultrametric Cluster Hierarchies: I Want 'em All!
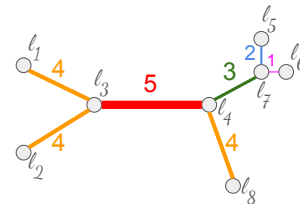
*Andrew Draganov\*, **Pascal Weber\***,
Rasmus Jørgensen, Anna Beer, Claudia Plant, Ira Assent*

December 3rd, 2025: 11 a.m. – 2 p.m.

Exhibit Hall C,D,E

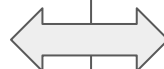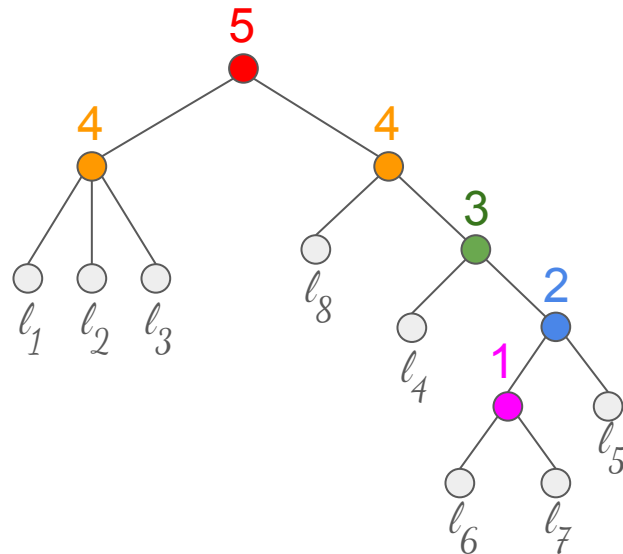# Different **Clustering** Algorithm Types



Centroid-based

Density-based

Graph-based

Hierarchical

Distribution-based

Ultrametric Cluster Hierarchies

# **All** ultrametrics are hierarchical

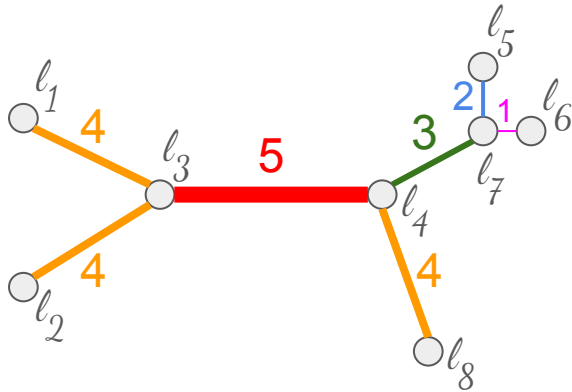

$$d(x, z) \leq \max \{d(x, y), d(y, z)\}$$



Distance between two leaves is the value in their lowest common ancestor (LCA):
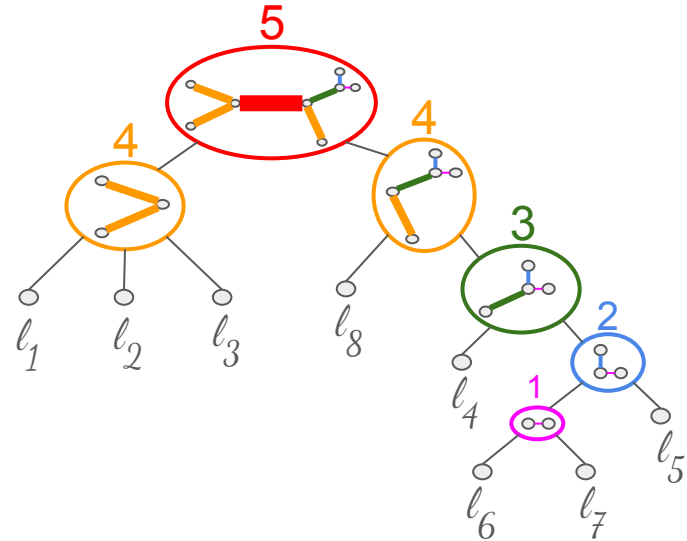
# **All** ultrametrics are hierarchical

Minimax distance is the largest step in the minimax path between two nodes, which is always on a minimum spanning tree (MST)

Distance between two leaves is the value in their lowest common ancestor (LCA):
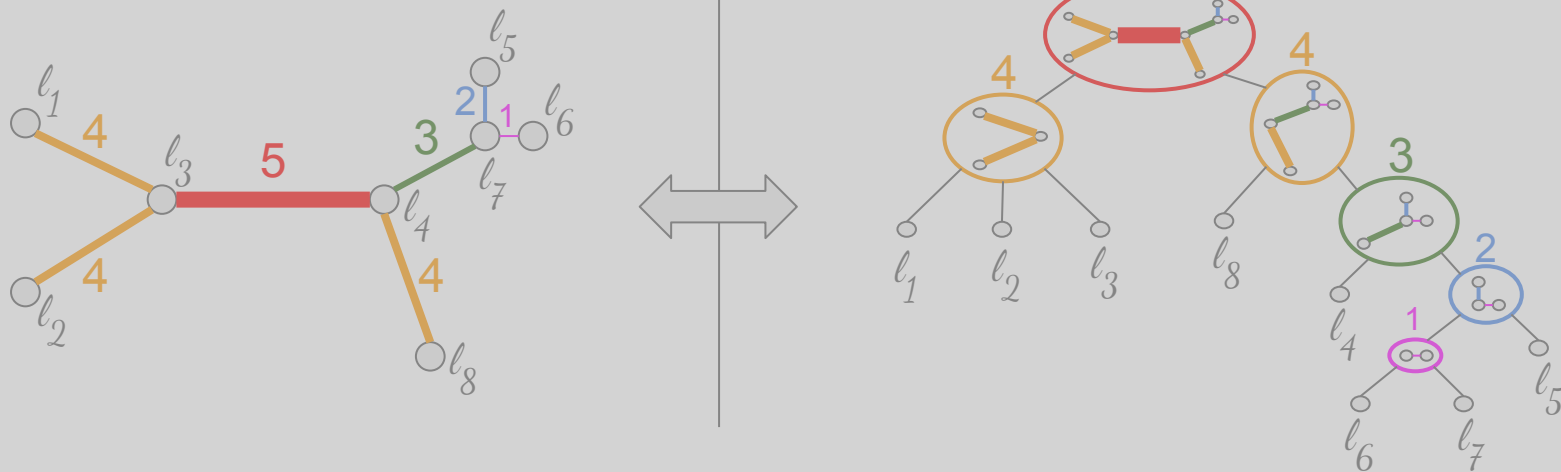
# **All** ultrametrics are hierarchical

Minimax distance is ~~the~~ ... ~~es~~ is the value
minimax path betwee~~n~~ ... ~~a~~ncestor (LCA):
always on a minimum spanning tree (MST)

**Density-based clustering**
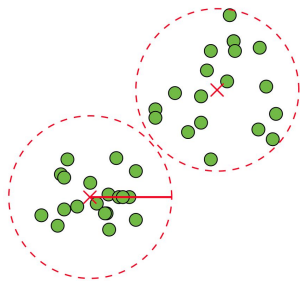corresponds to thresholding an **ultrametric**

# What happens if we do centroid-based clustering in ultrametrics?

- The following are NP-hard in standard metric spaces
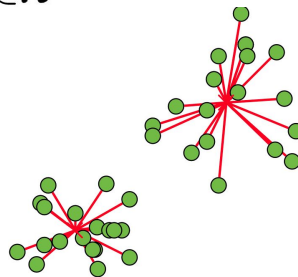  but can be **solved optimally in ultrametrics**:

$k$-center

$$\max_{x \in \mathcal{X}} \min_{c \in C} d(x, c)$$

$k$-means

$$\sum_{x \in \mathcal{X}} \min_{c \in C} d(x, c)^2$$

# What happens if we do centroid-based clustering in ultrametrics?

- The follow[...] It takes **`Sort(n) time`** to find the optimal $k$-means
  but can be solutions for **all** values of $k$ **in an ultrametric**

$k$-center

$$\max_{x \in \mathcal{X}} \min_{c \in C} d(x, c)$$

$k$-means

$$\sum_{x \in \mathcal{X}} \min_{c \in C} d(x, c)^2$$

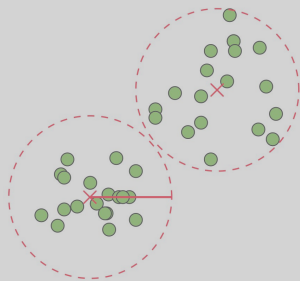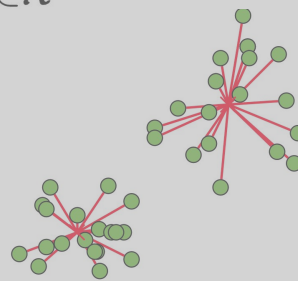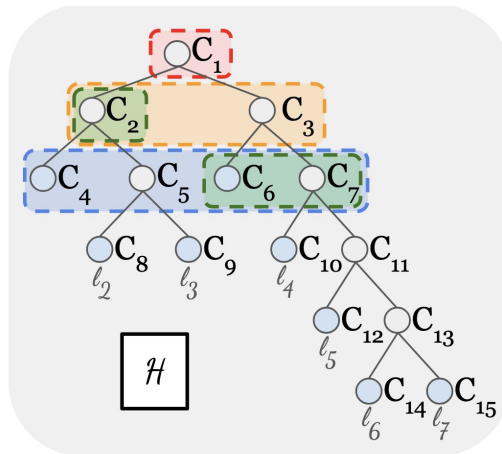# The optimal *k*-means solutions are *themselves* hierarchical

The optimal *k*-means solution and the optimal (*k*+1)-means solution in an ultrametric are <u>identical</u> except that a single cluster was split in two
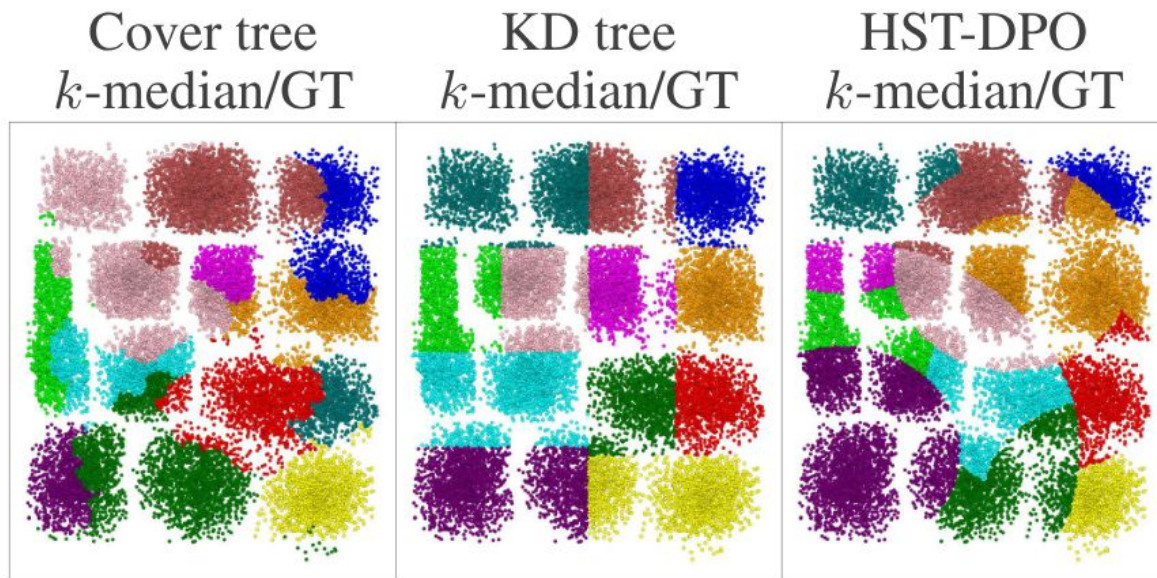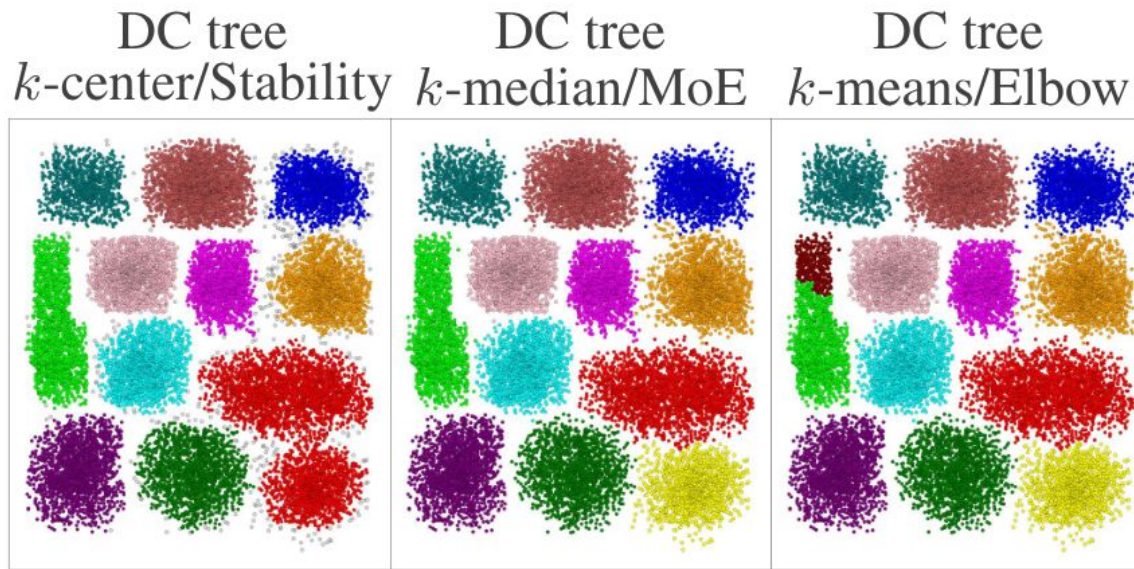
# Implications

1. Solve *any* center-based clustering task in *any* ultrametric *optimally* in `Sort(n)` time
   a. Takes this time for *all* values of *k*
   b. These solutions are hierarchical

2. One can pick any clustering from this hierarchy extremely quickly (in $O(n)$ time)
   a. Threshold the values in the tree (DBSCAN)
   b. Pick the "best" clustering by a function (HDBSCAN)
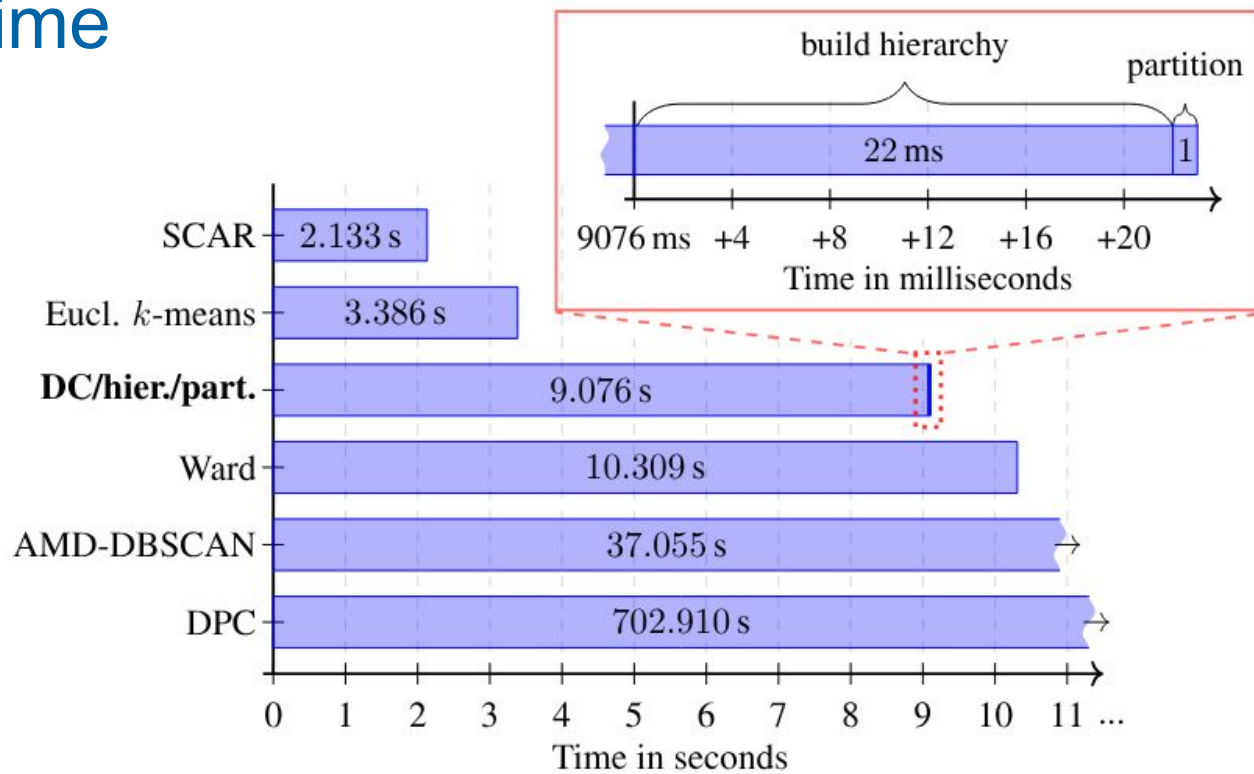   c. Optimal clustering for a user-specified value of *k*
   d. Elbow method

# Different Ultrametrics

# Different Hierarchies

# Runtime

# Summary

- It takes **Sort(n) time** to find the optimal $k$-means solutions for **all** values of $k$ **in an ultrametric**

- Outperforms other clustering algorithms in clustering accuracy

- Ultrametric Cluster Hierarchies – Implementation:
  **S**imilarity **Hi**erarchy **P**artitioning Framework (SHiP framework)

Github Repo

pip package