# Dask Summit 2021
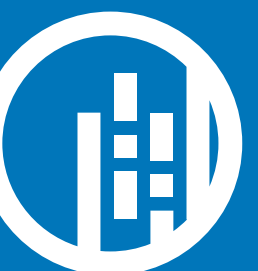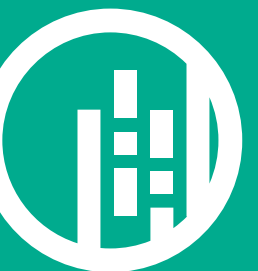
# Scaling geospatial vector data

## Partitioning of spatial data

**Martin Fleischmann**
@martinfleis

# How Dask chunks data?

NumPy Array

Dask Array

**dask.array**

January, 2016

February, 2016 } Pandas Dataframe

March, 2016 } Dask Dataframe

April, 2016

May, 2016

+

NumPy Array

Dask Array
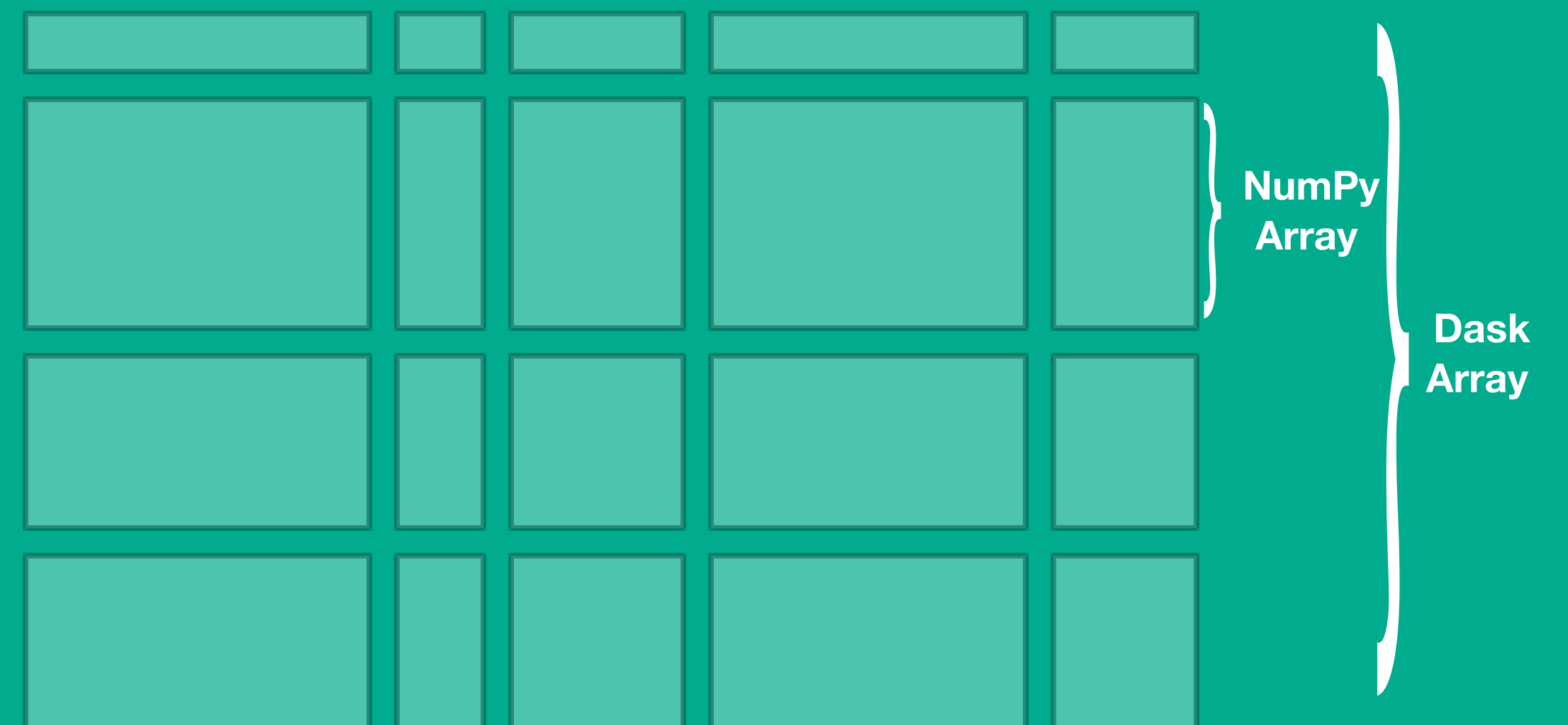
spatial dataframe

# Why we need spatially coherent chunks?

**spatial indexing**

**spatial operations**

**spatial operations**

spatial proximity
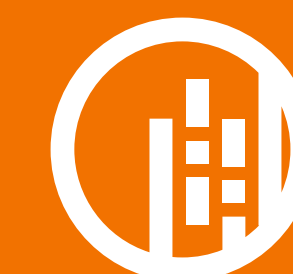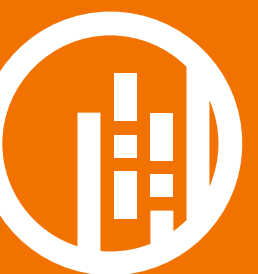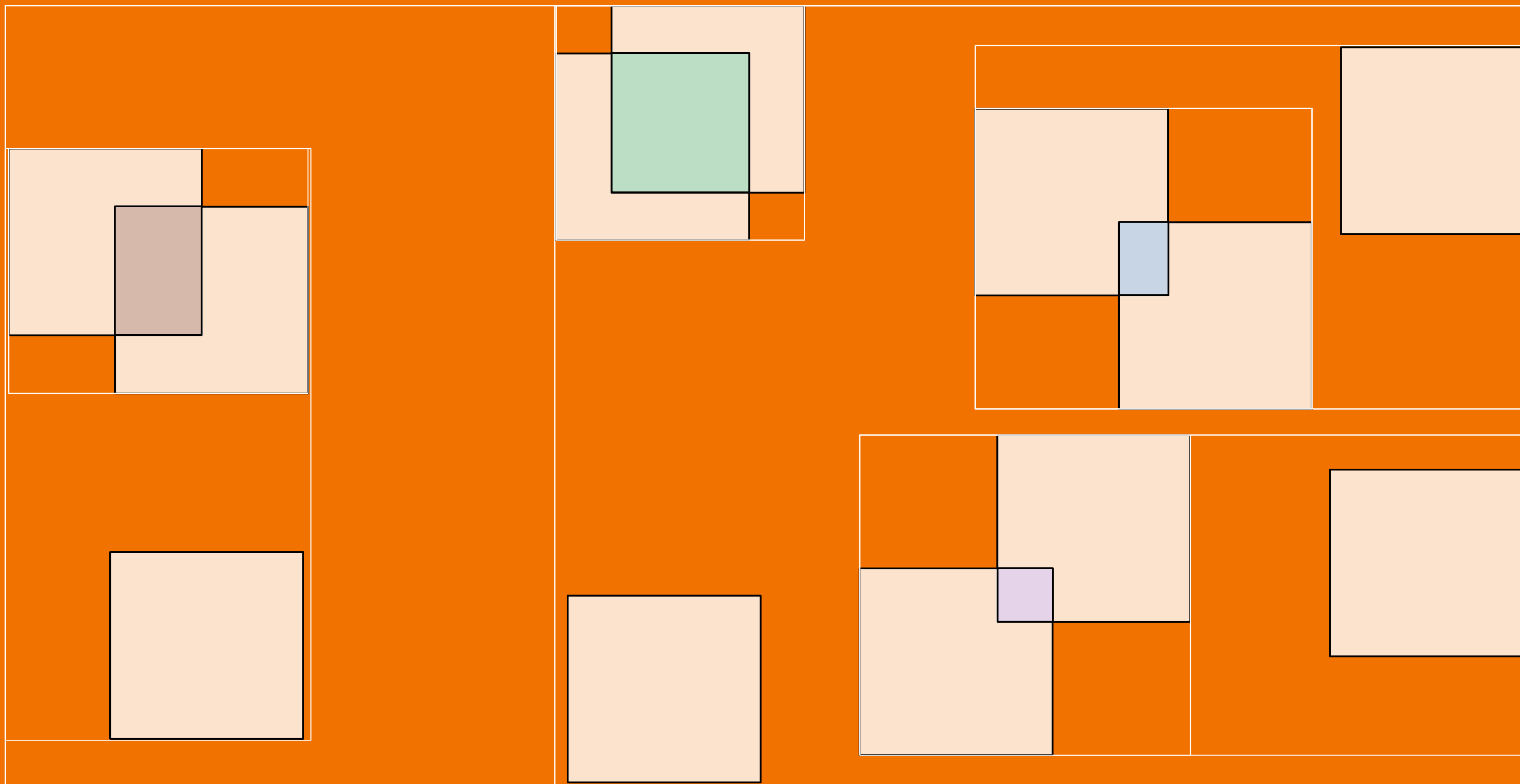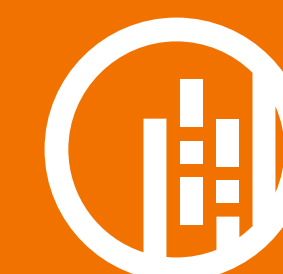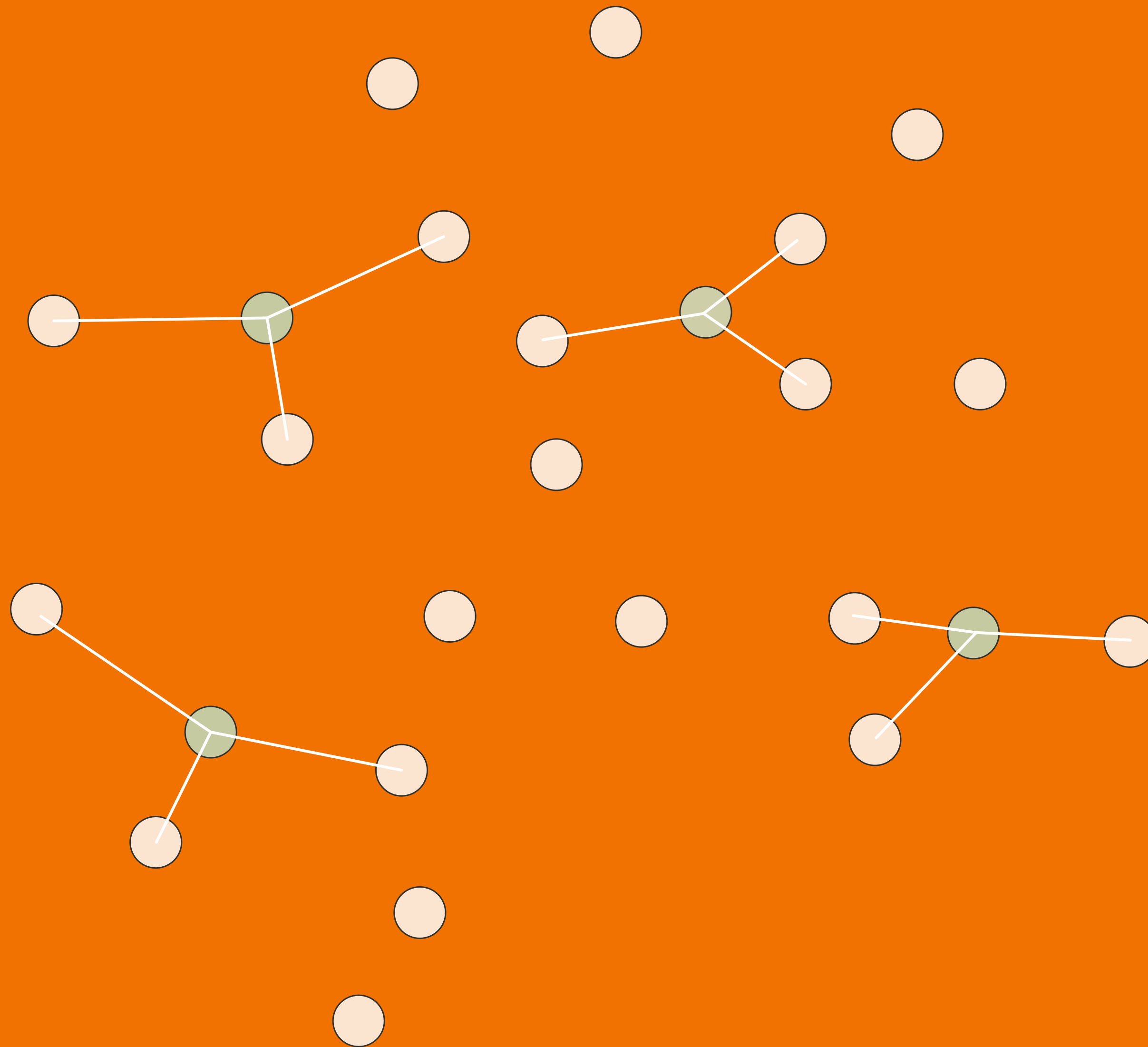
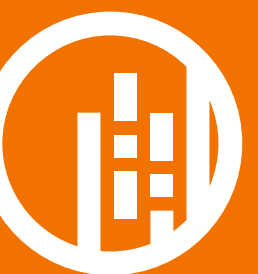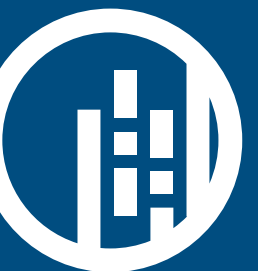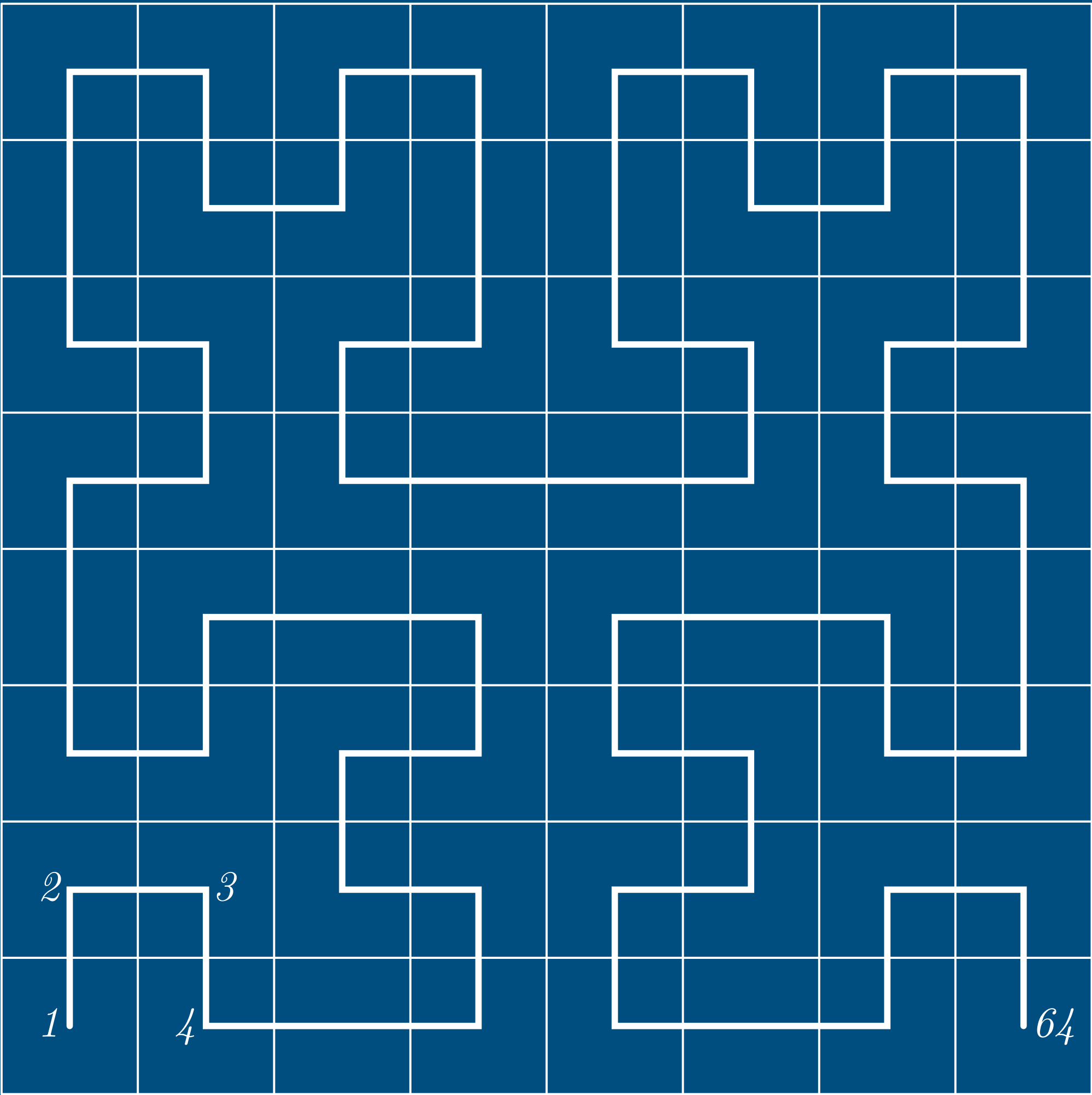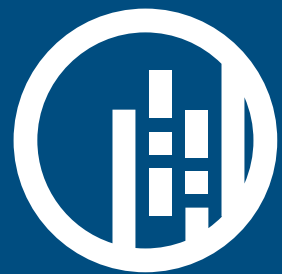**Pandas Dataframe**

**Dask Dataframe**
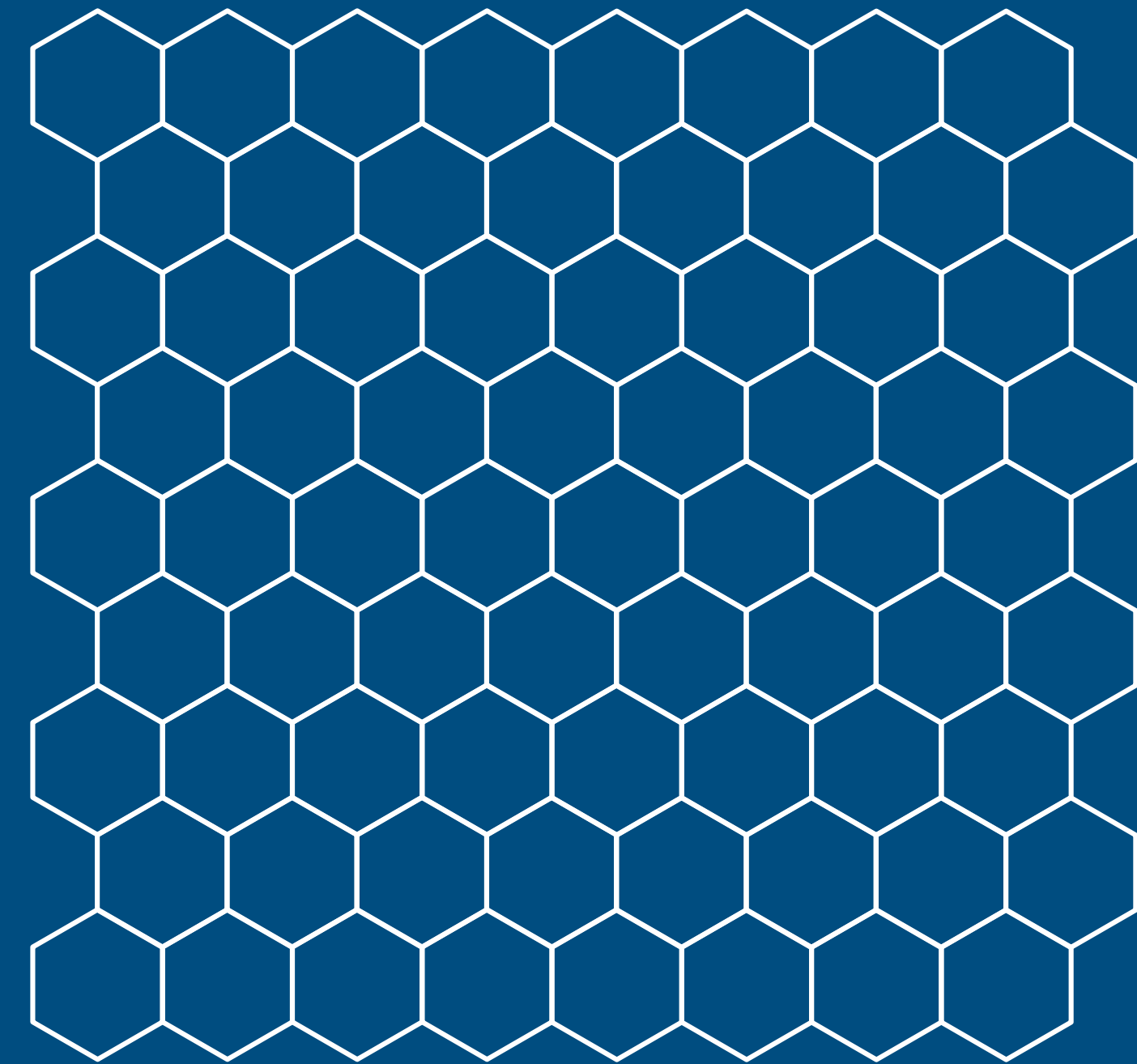
**NumPy Array**

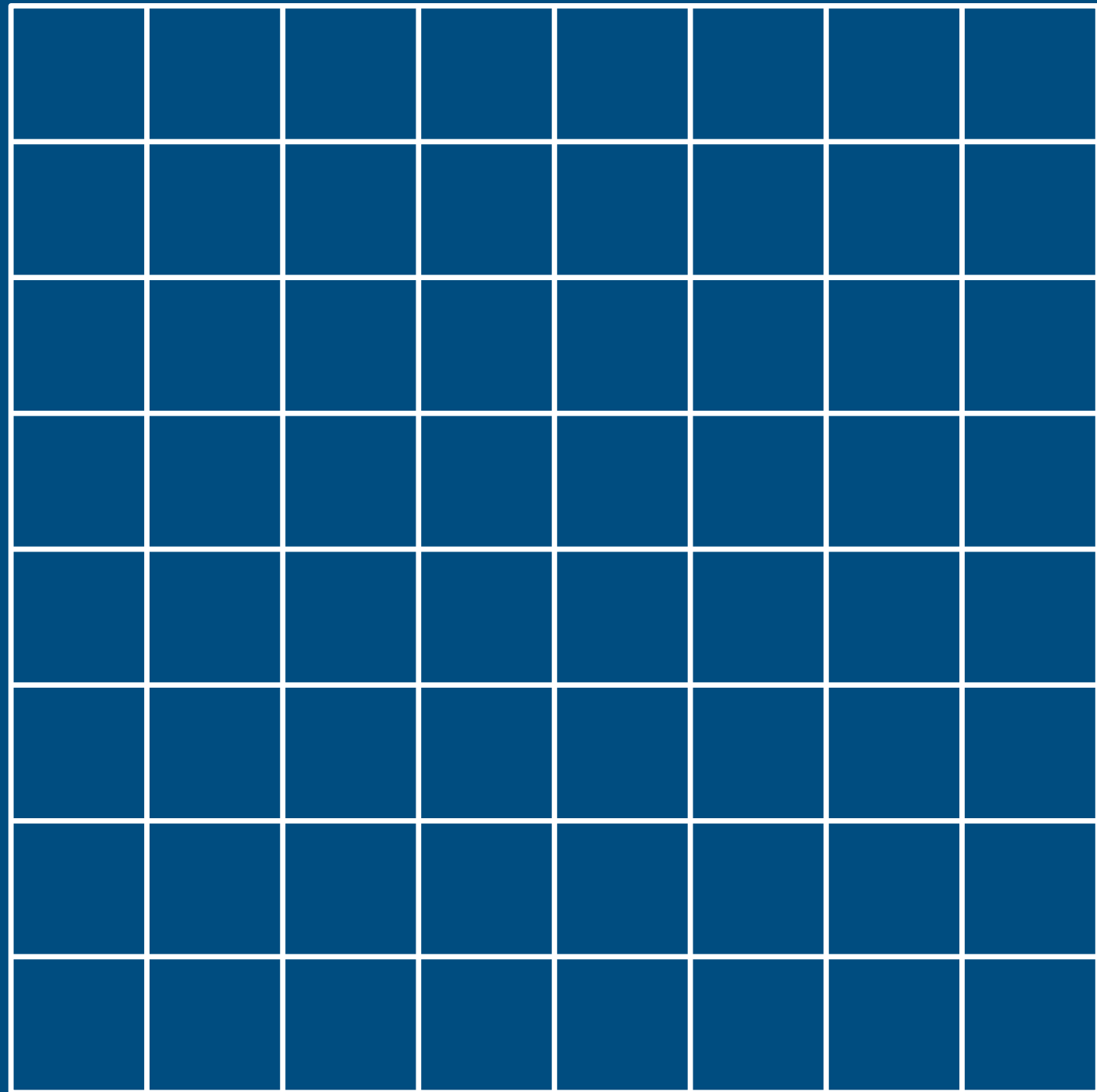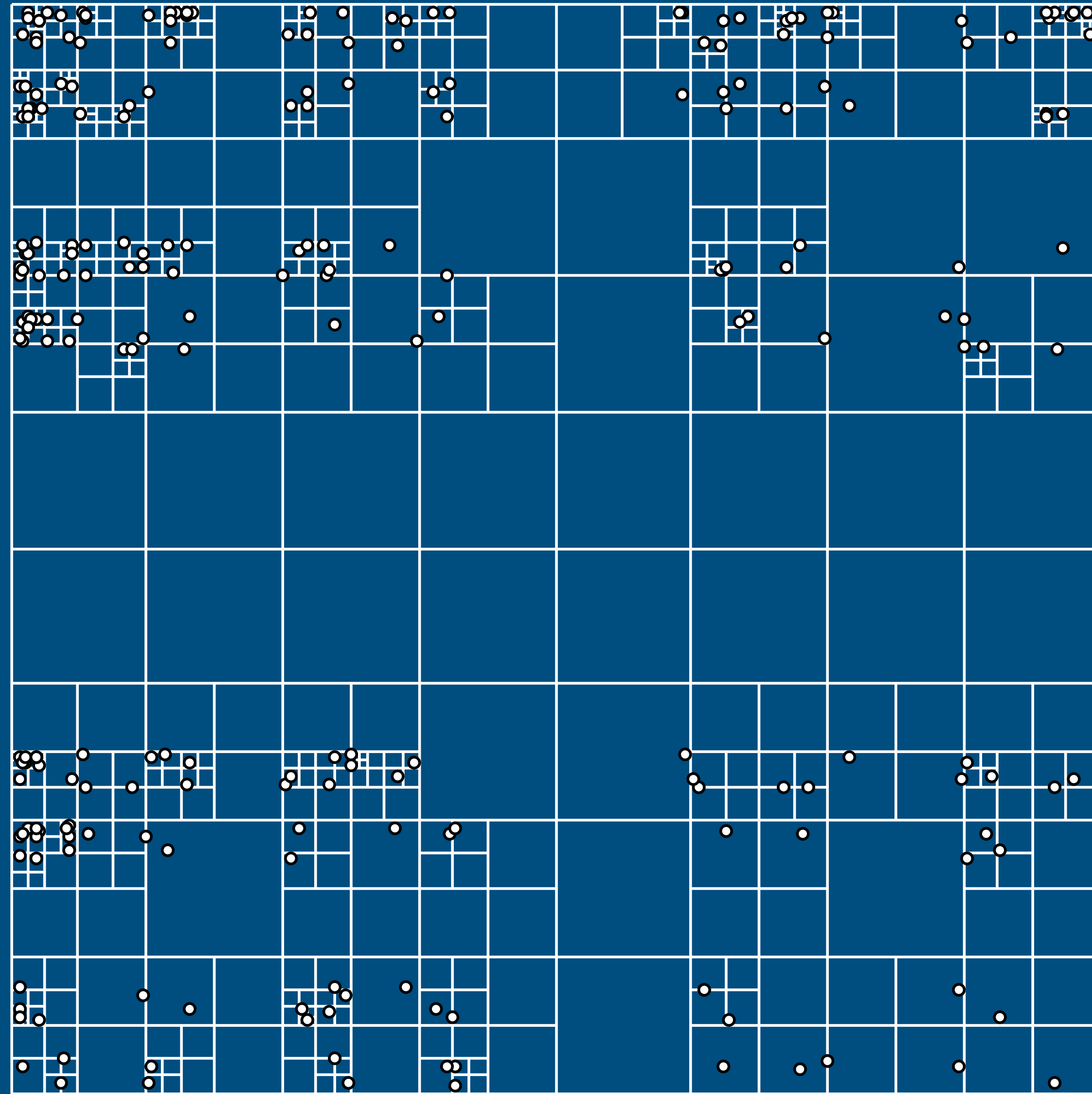**Dask Array**

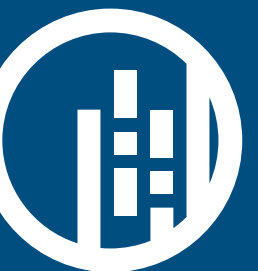# overlapping computation
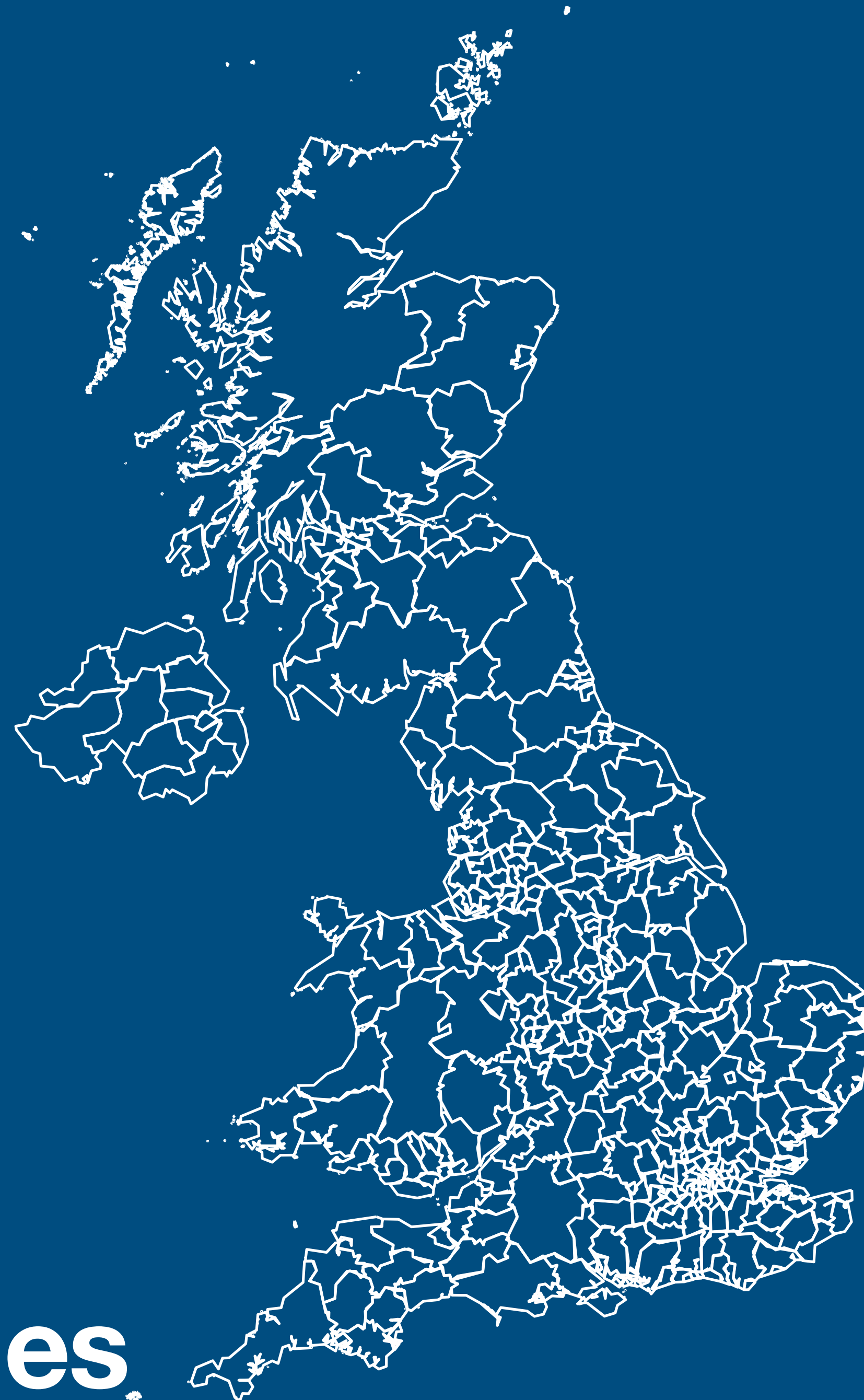
# How can we achieve that?

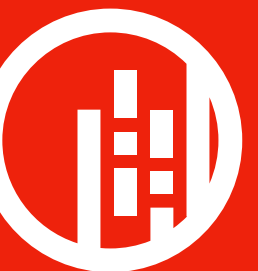**Hilbert curve**

**(arbitrary) grid**

quadtree

known boundaries

# When it gets tricky

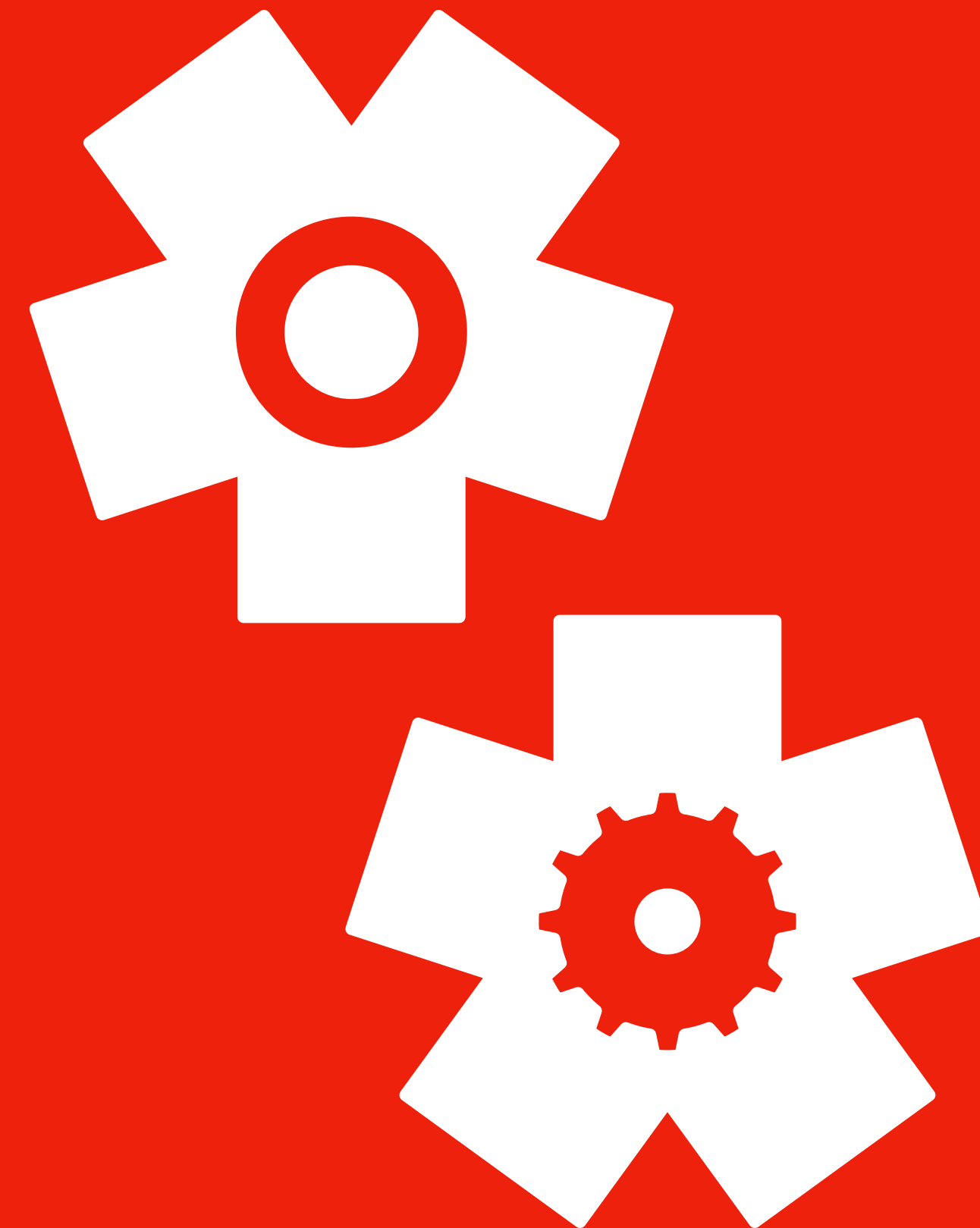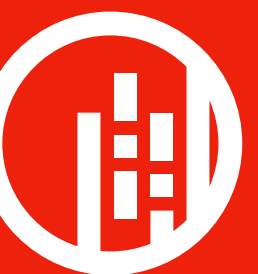# vector data are unpredictable
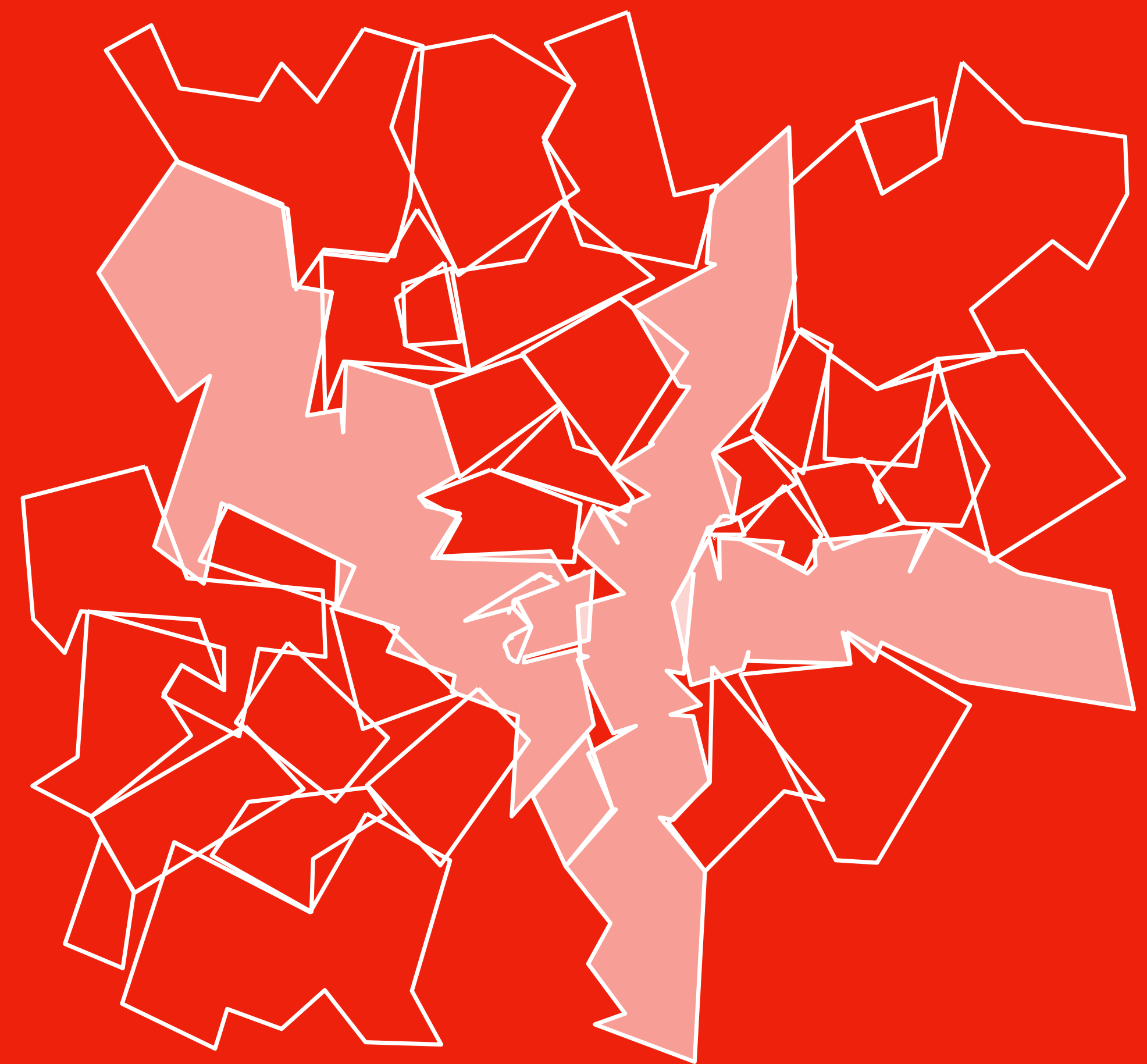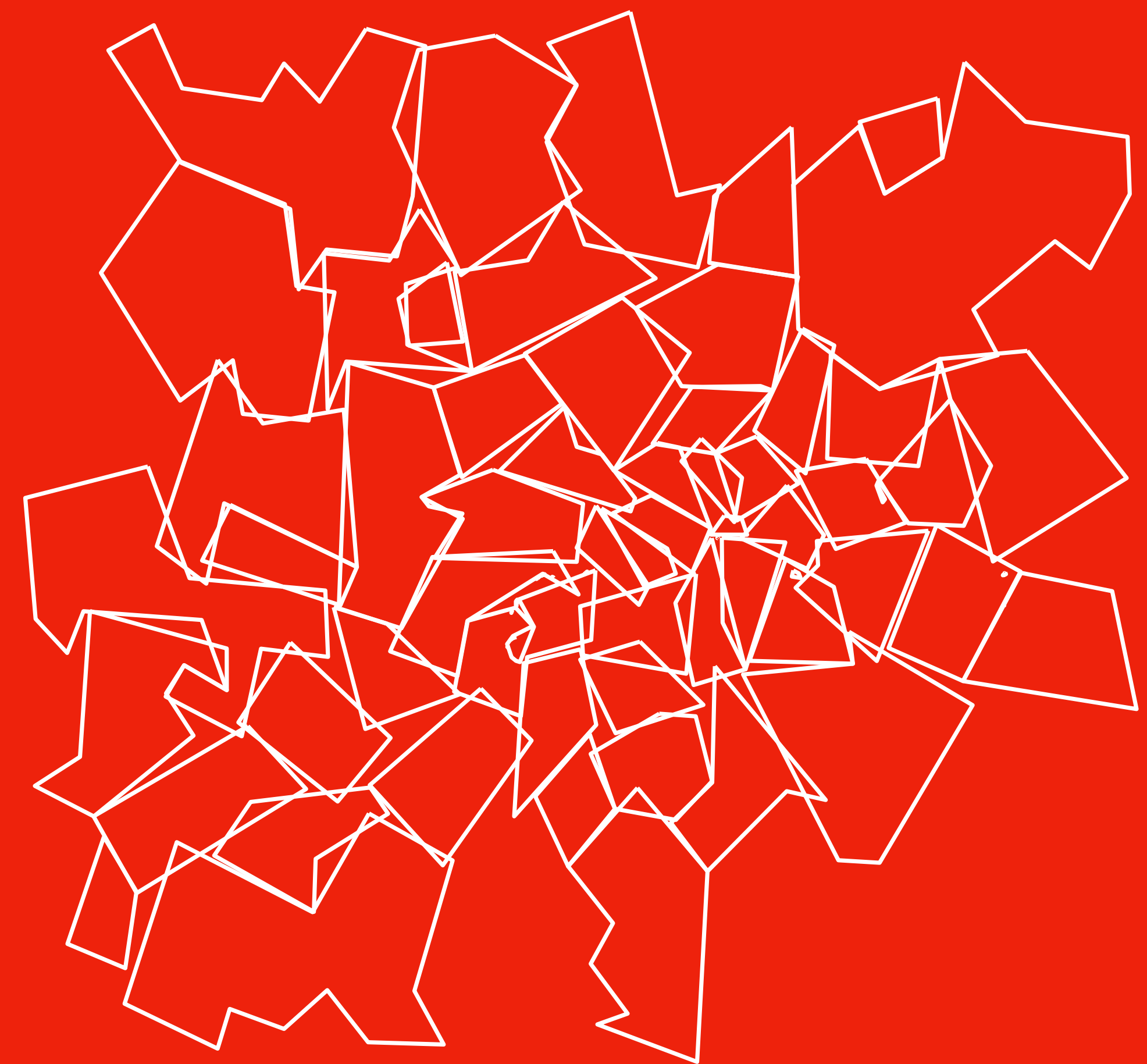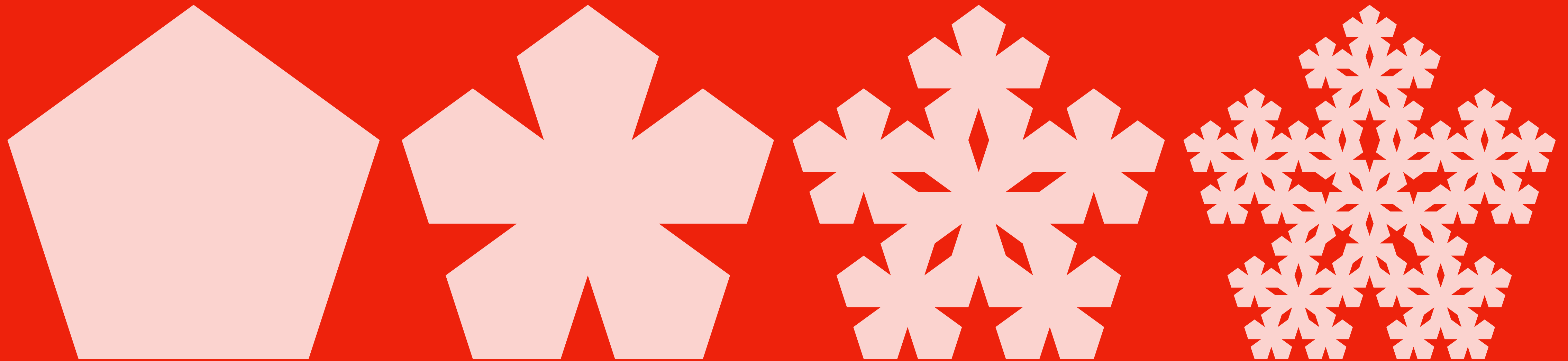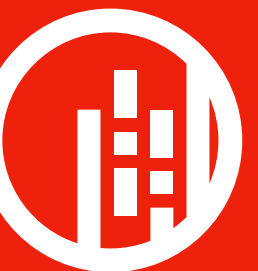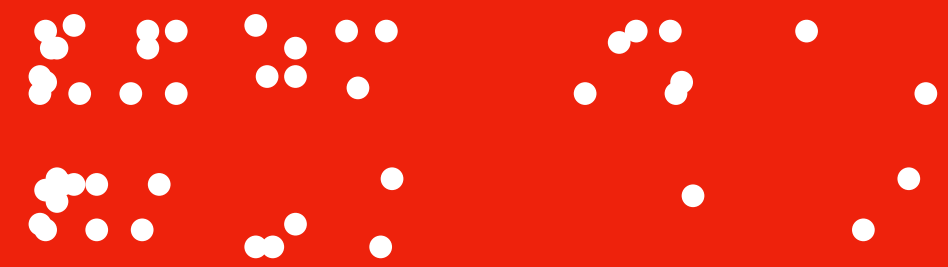
geometry types

geometry types

**boundaries**

boundaries

complexity
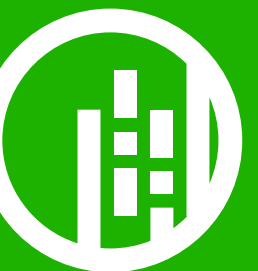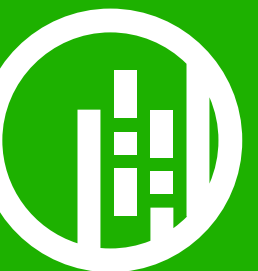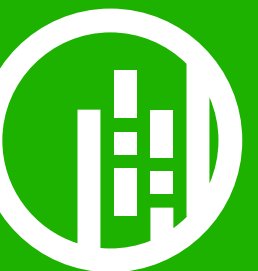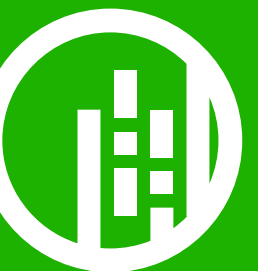
spatial distribution

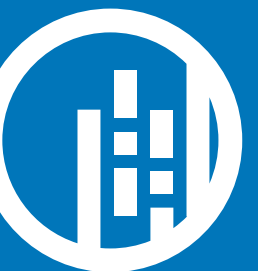# What would be nice to have

# efficient cross-chunk spatial indexing

support for overlapping computations

# Take into account memory demands of different rows?

# Questions

# How to create *smart* spatial partitions?

# How to store spatial partitions?