# Aggregation of Dots

## Methods for Big Data in Official Statistics

Martijn Tennekes

Heerlen, October 5, 2018

*132,735,324 dots in this presentation!*
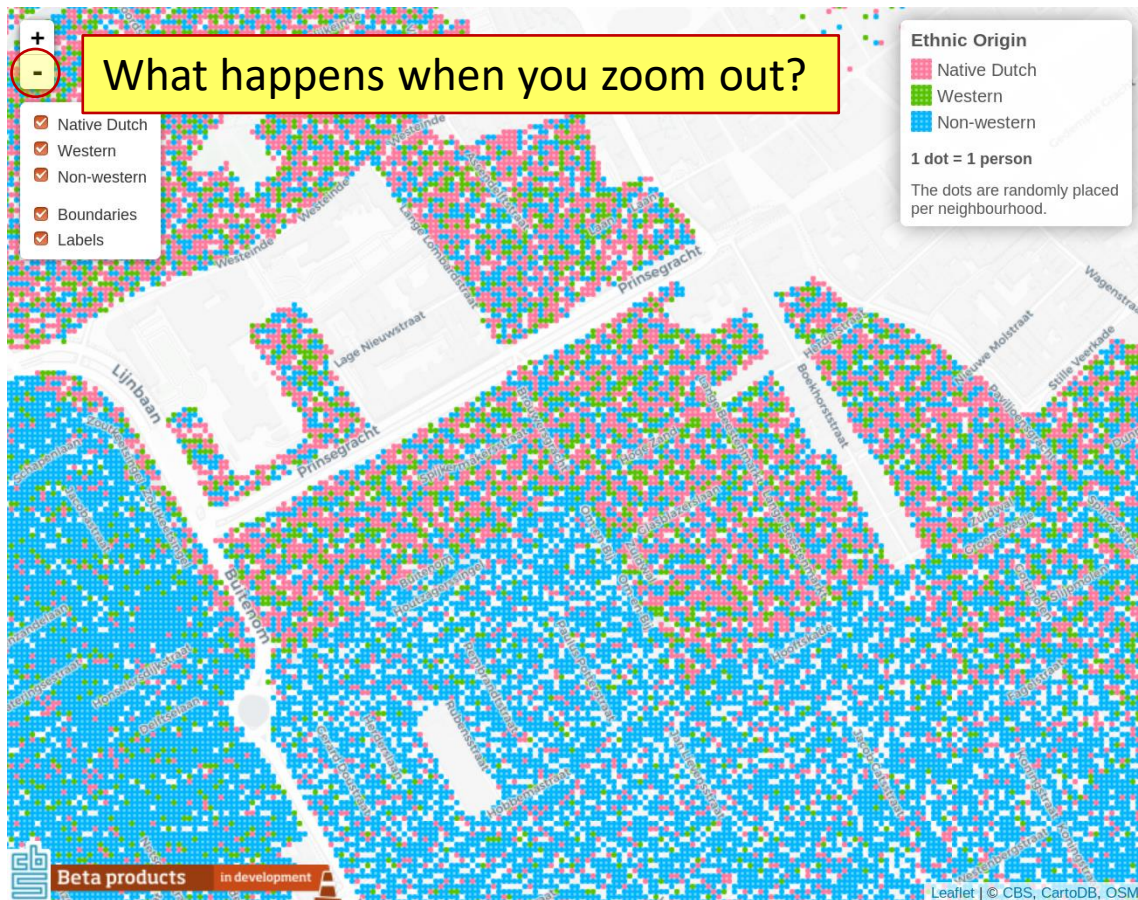
# Classic dot map



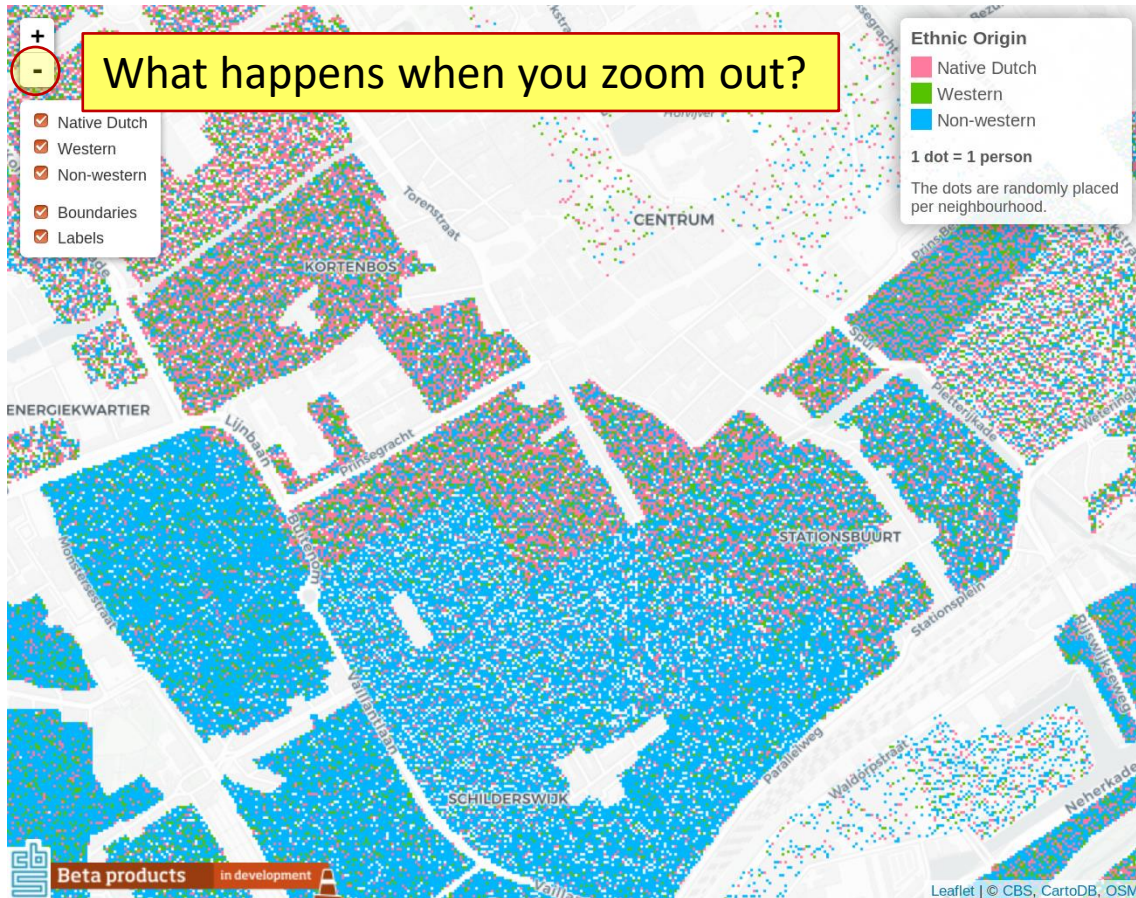Cholara outbreak in London (1854) by John Snow    Dots instead of bars    2

# Let there be… COLOR



Position of the dots:
**density**

Colors of the dots:
**composition**

3

# What happens when you zoom out?



Position of the dots:
**density**

Colors of the dots:
**composition**

# Out of pixels ☹

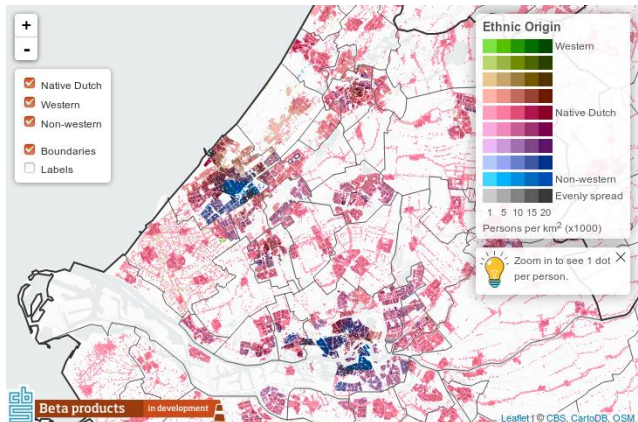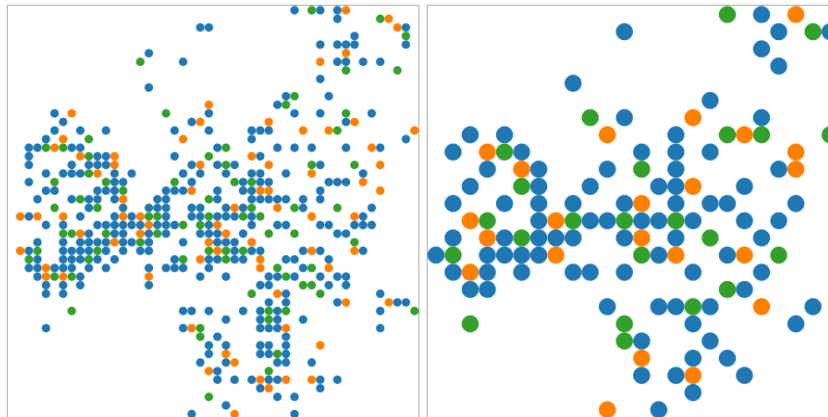Wait until 8K UHD becomes the standard?     Nope.

*Hmm, still can't see the dots...*

How to aggregate the dots?

We propose two approaches:
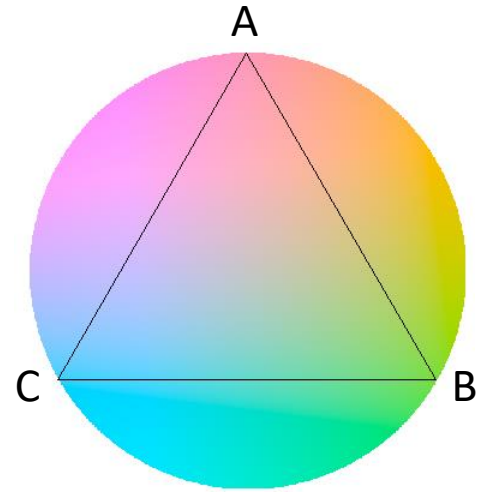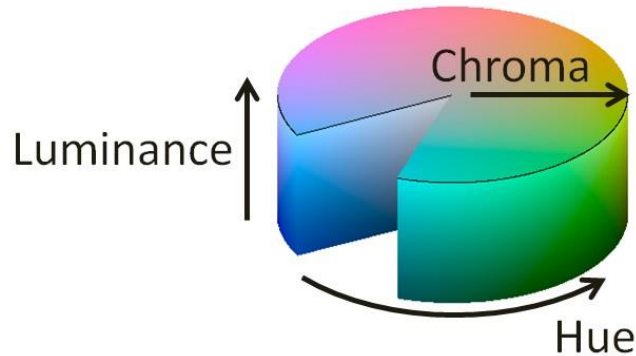
1. Blended colours

2. Super dots

# Blended colours

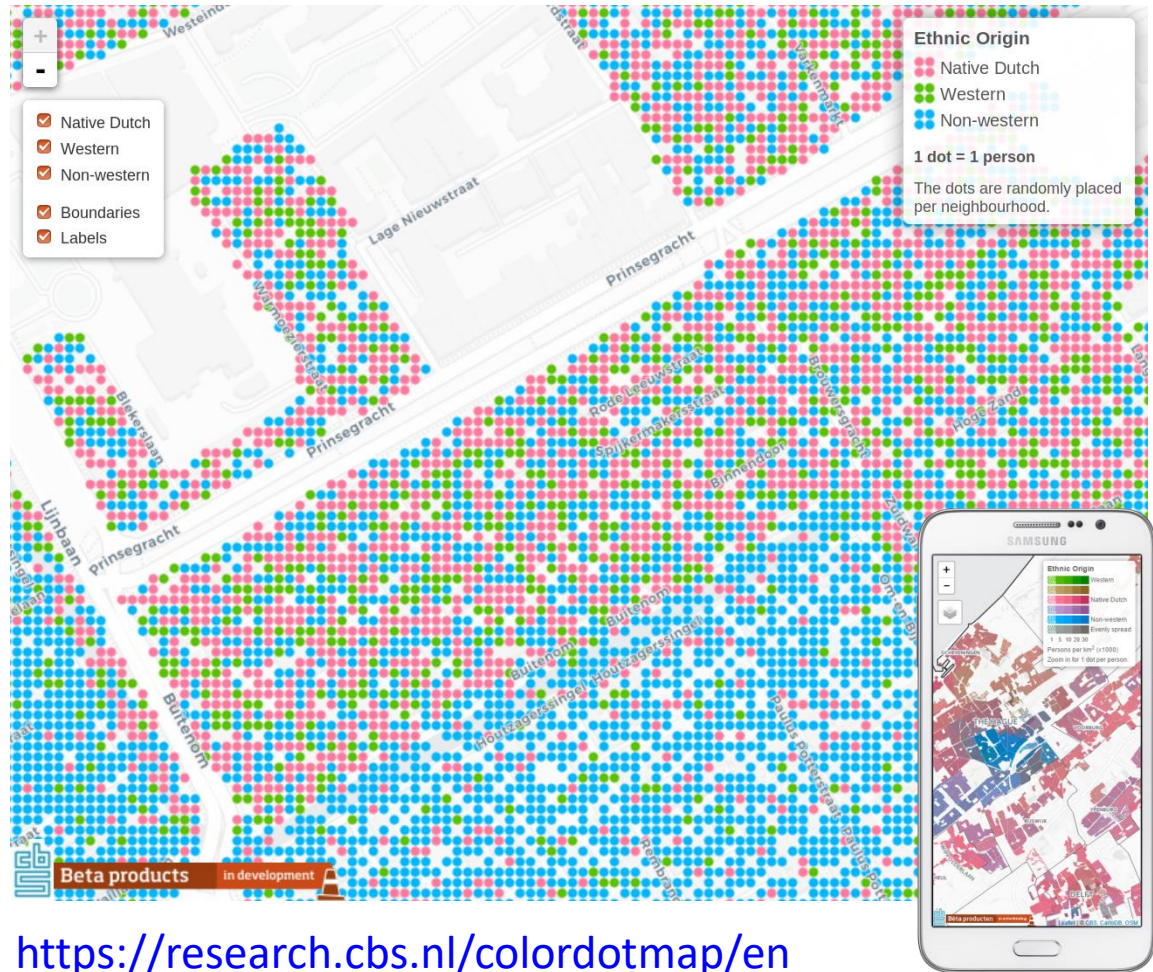Pixel colours are selected from the HCL colour space:





- **Luminance** for **density**
- **Hue** and **Chroma** for **composition**

# **Application**

Migration background of the Dutch population

Dots are distributed uniformly per neighbourhood and placed in the land use category "residential"



Published a CBDS beta product: https://research.cbs.nl/colordotmap/en

# Application

Migration background of the Dutch population

Dots are distributed uniformly per neighbourhood and placed in the land use category "residential"



Published a CBDS beta product: https://research.cbs.nl/colordotmap/en

# **Application**

Migration background of the Dutch population

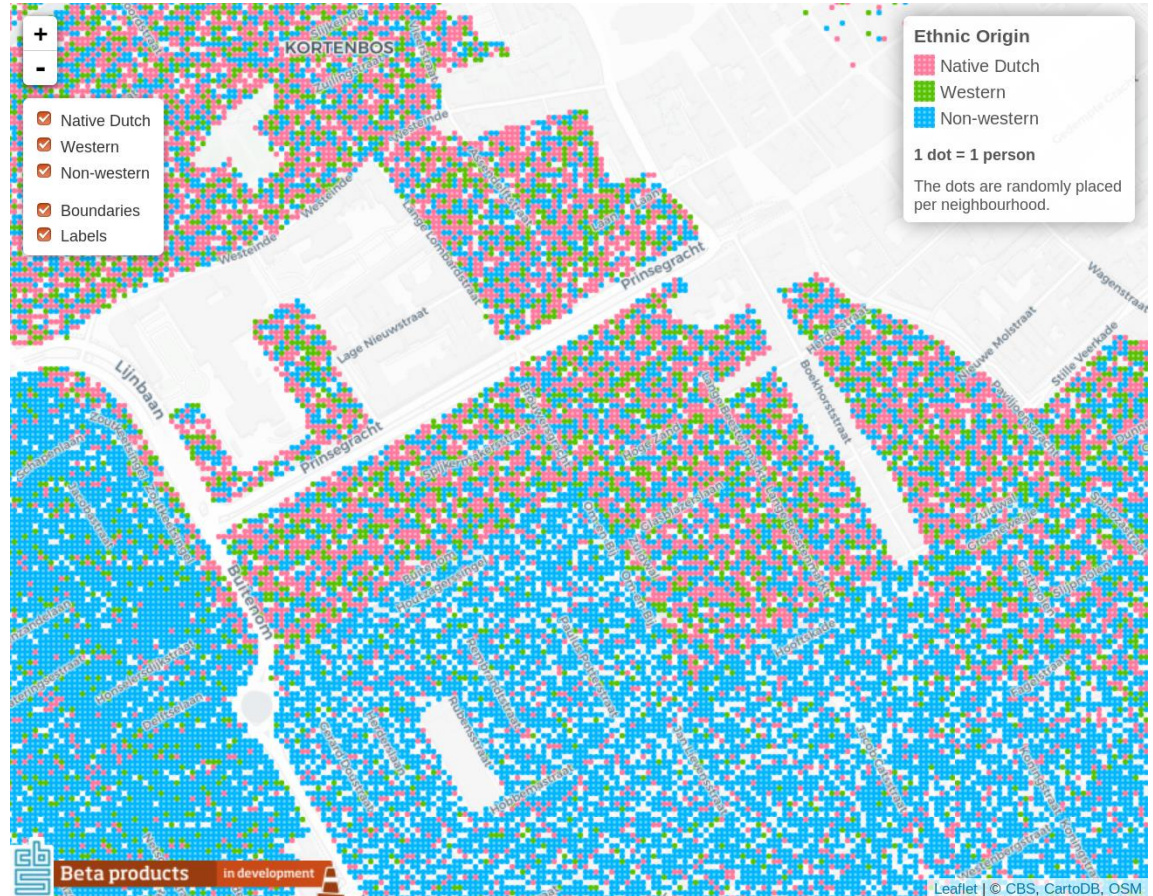Dots are distributed uniformly per neighbourhood and placed in the land use category "residential"



Published a CBDS beta product: https://research.cbs.nl/colordotmap/en

# Application

Migration background of the Dutch population

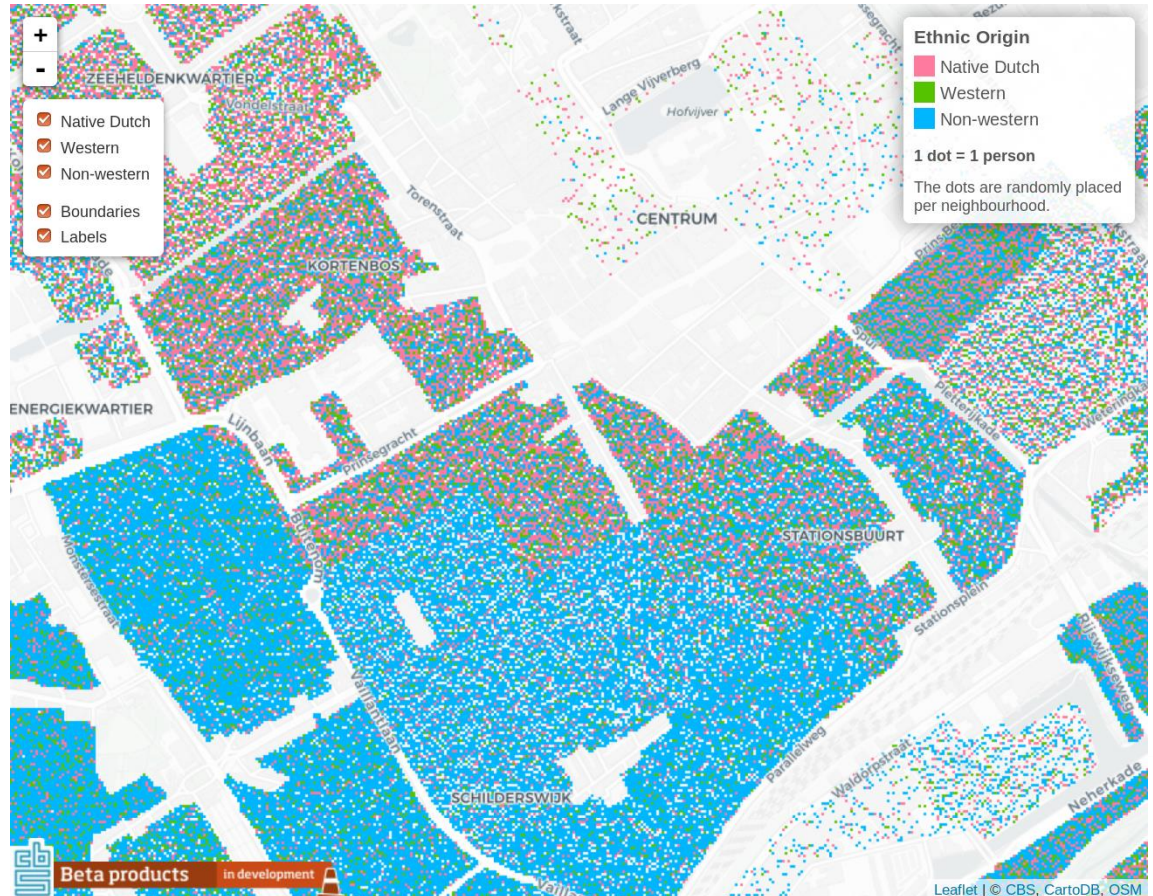Dots are distributed uniformly per neighbourhood and placed in the land use category "residential"



Published a CBDS beta product: https://research.cbs.nl/colordotmap/en

# Application

Migration background of the Dutch population

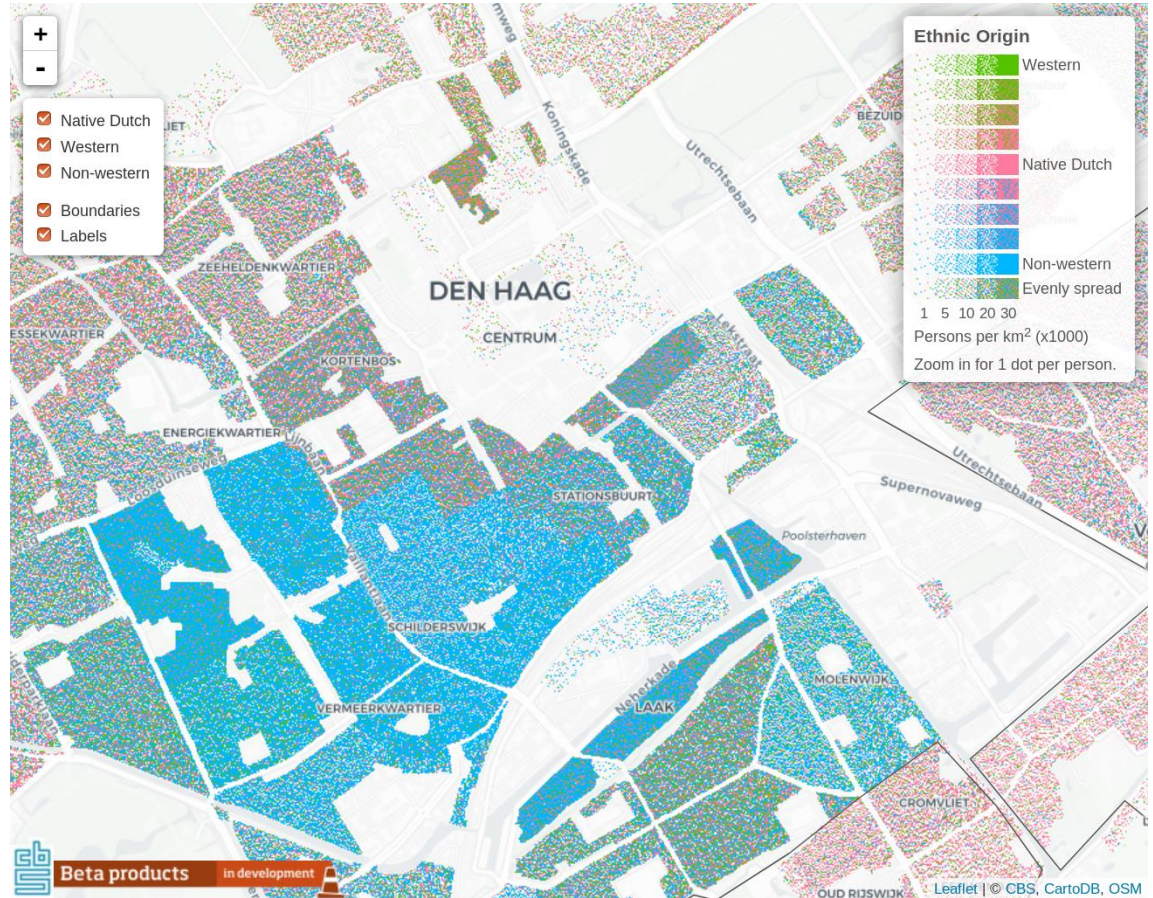Dots are distributed uniformly per neighbourhood and placed in the land use category "residential"



Published a CBDS beta product: https://research.cbs.nl/colordotmap/en

# Application

Migration background of the Dutch population

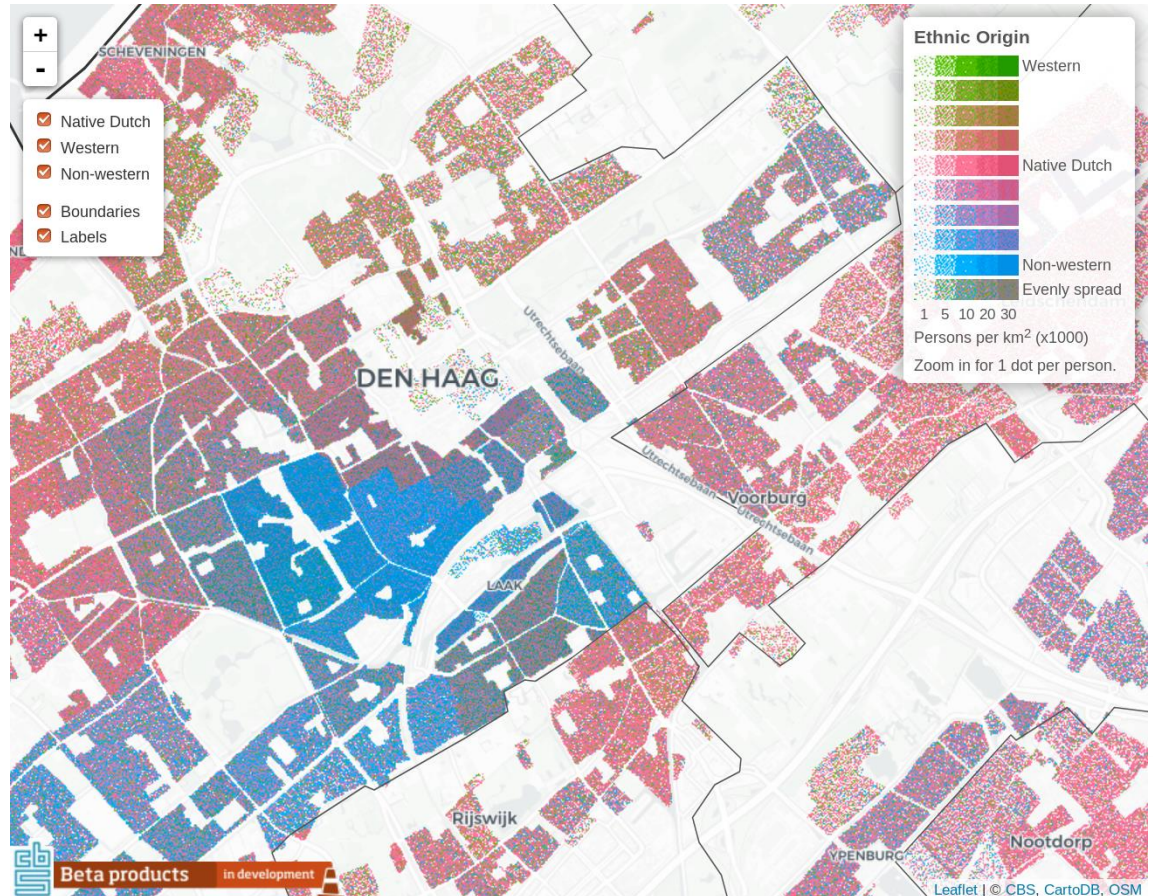Dots are distributed uniformly per neighbourhood and placed in the land use category "residential"



Published a CBDS beta product: https://research.cbs.nl/colordotmap/en

# Application

Migration background of the Dutch population

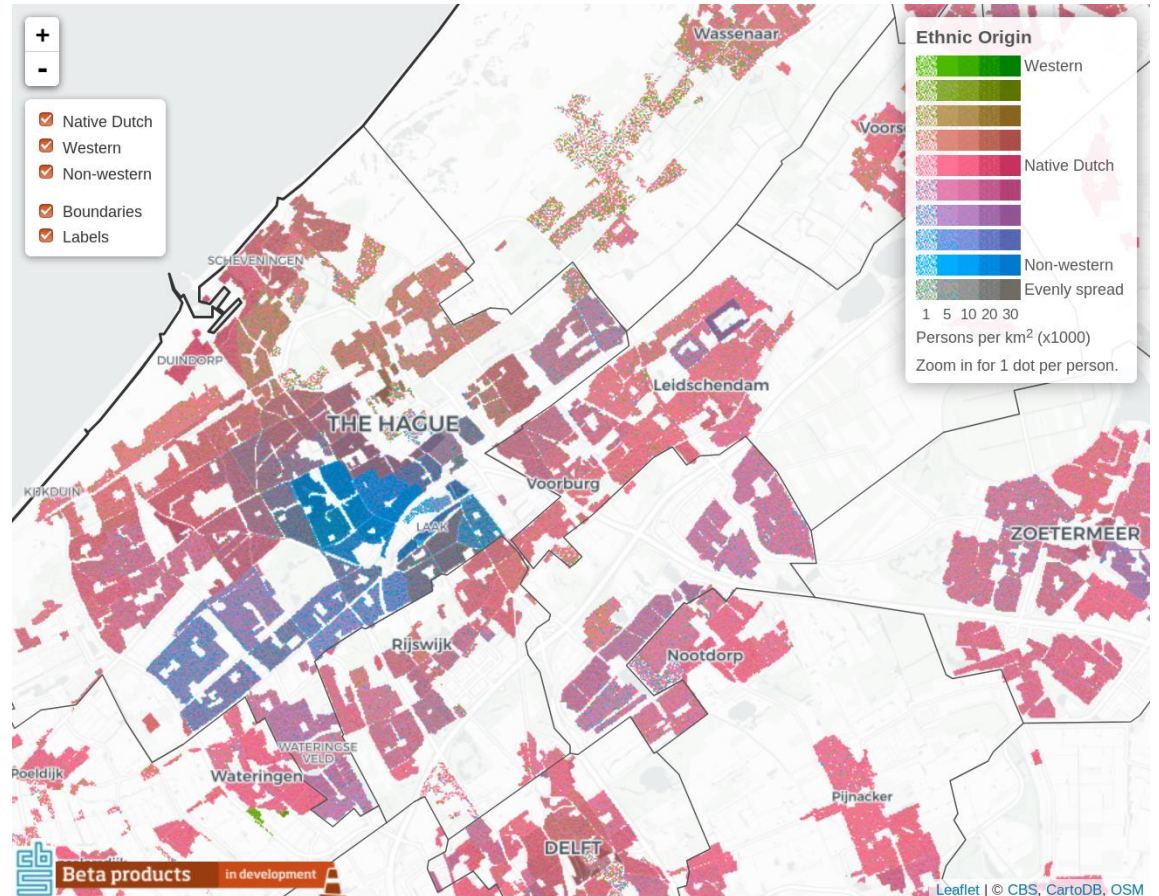Dots are distributed uniformly per neighbourhood and placed in the land use category "residential"



Published a CBDS beta product: https://research.cbs.nl/colordotmap/en

# Application

Migration background of the Dutch population

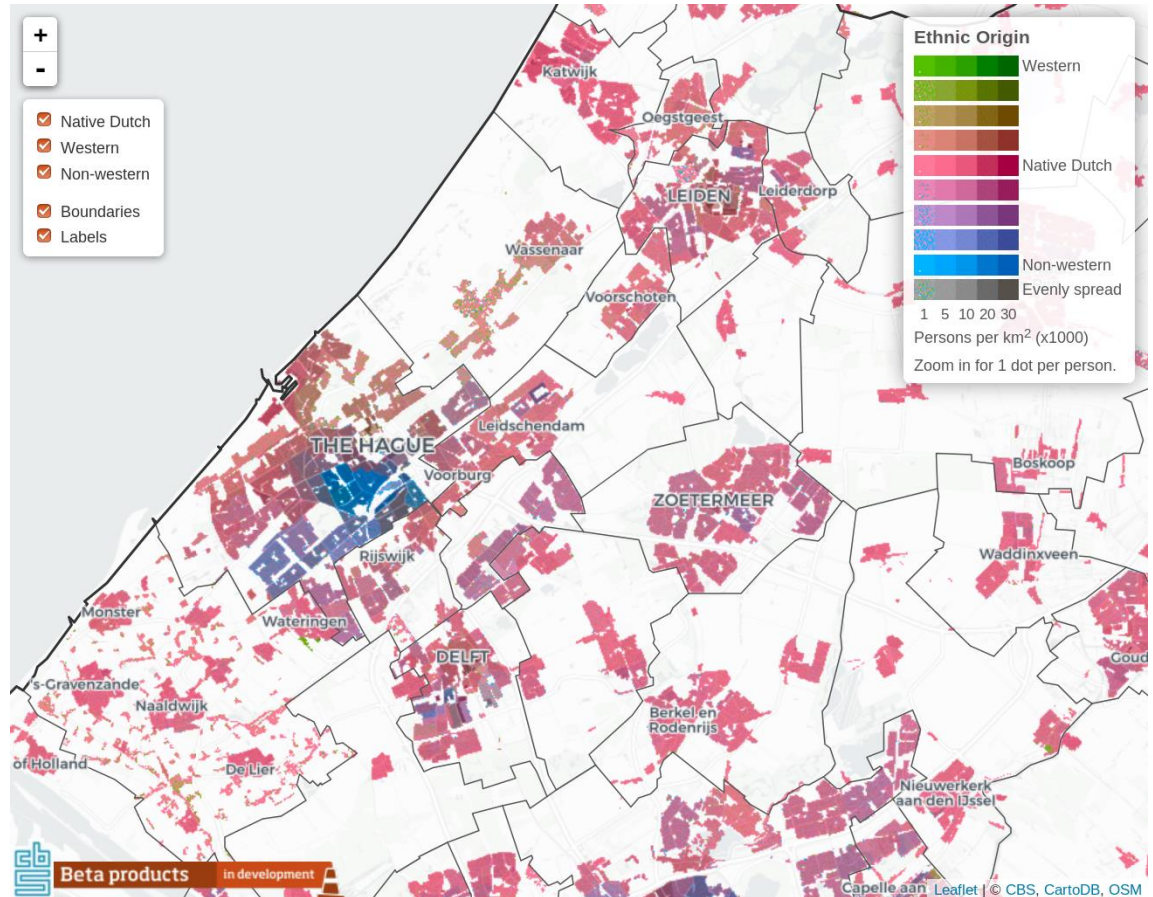Dots are distributed uniformly per neighbourhood and placed in the land use category "residential"



Published a CBDS beta product: https://research.cbs.nl/colordotmap/en

# Application

Migration background of the Dutch population

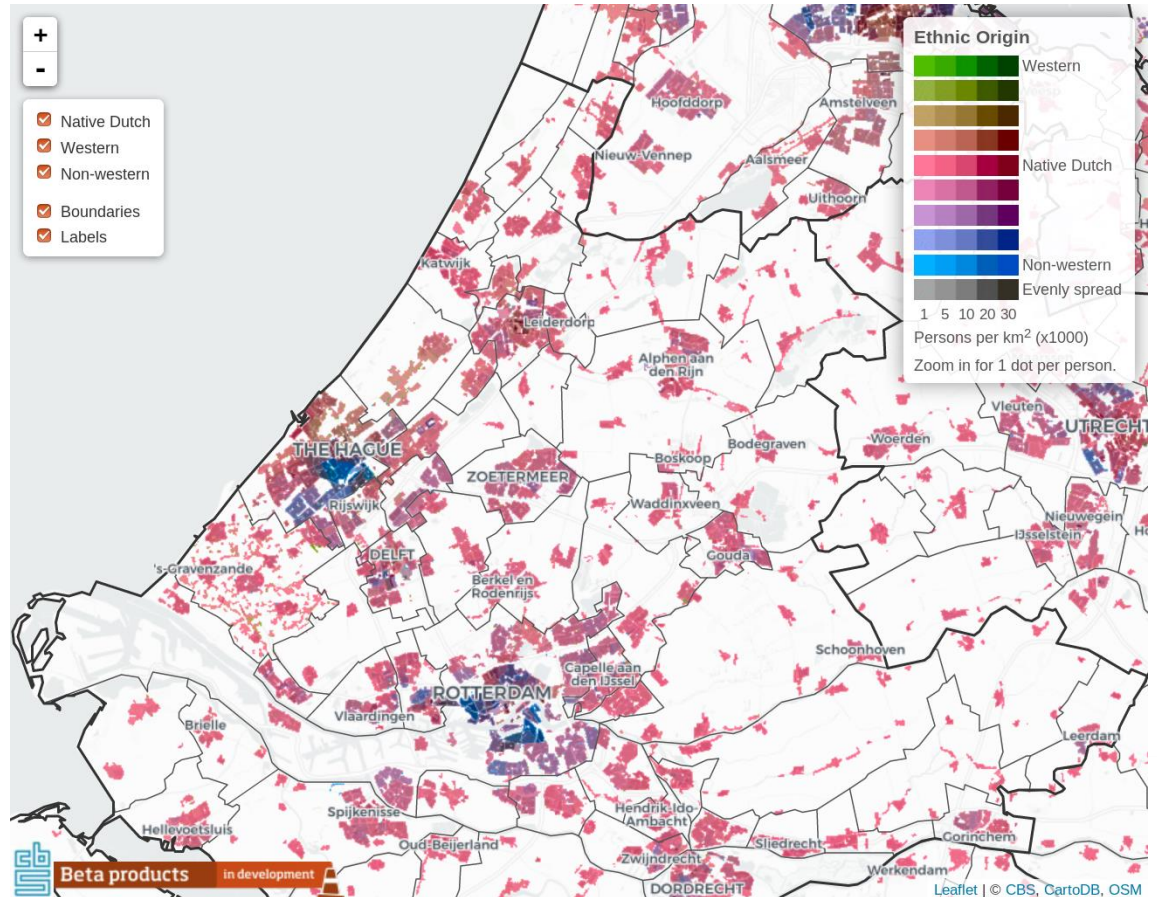Dots are distributed uniformly per neighbourhood and placed in the land use category "residential"



Published a CBDS beta product: https://research.cbs.nl/colordotmap/en

# Application

Migration background of the Dutch population

Dots are distributed uniformly per neighbourhood and placed in the land use category "residential"



Published a CBDS beta product: https://research.cbs.nl/colordotmap/en

# **Application**

Migration background of the Dutch population

Dots are distributed uniformly per neighbourhood and placed in the land use category "residential"
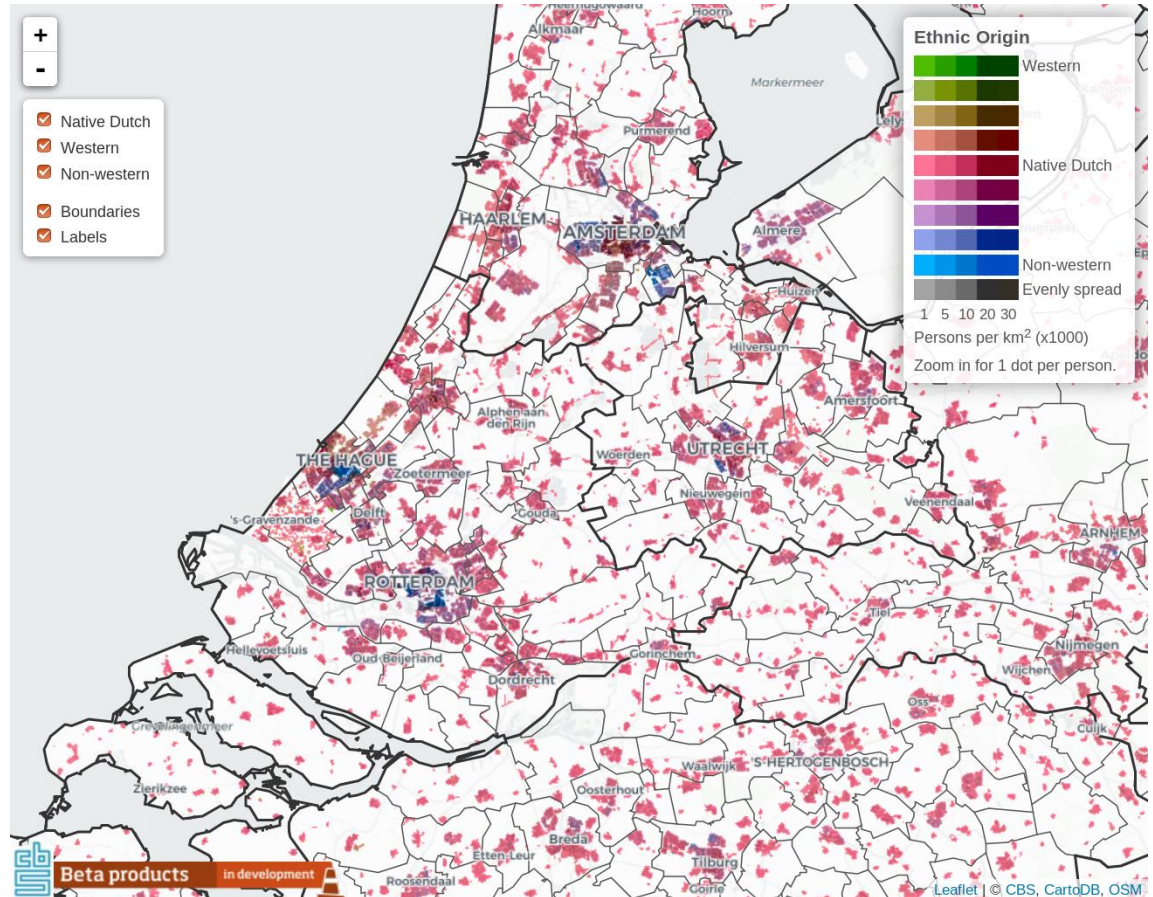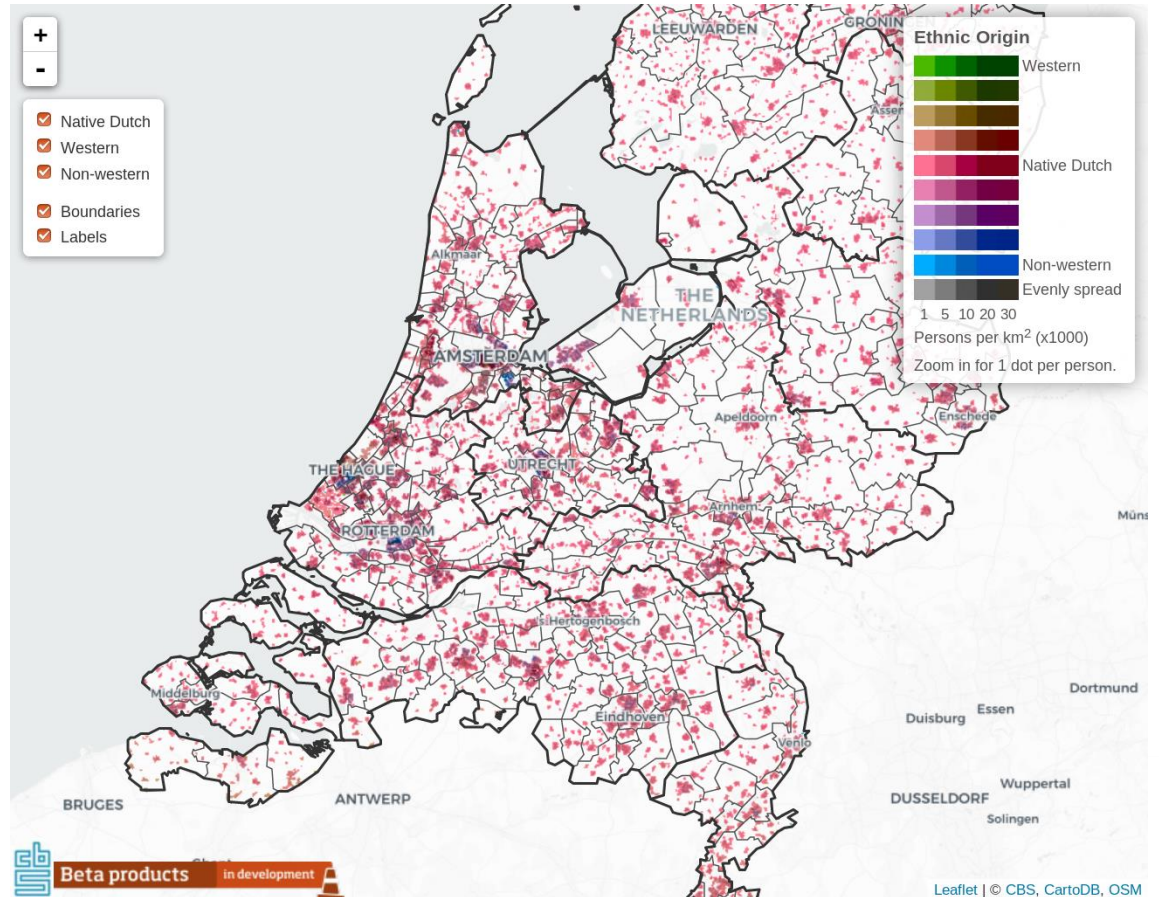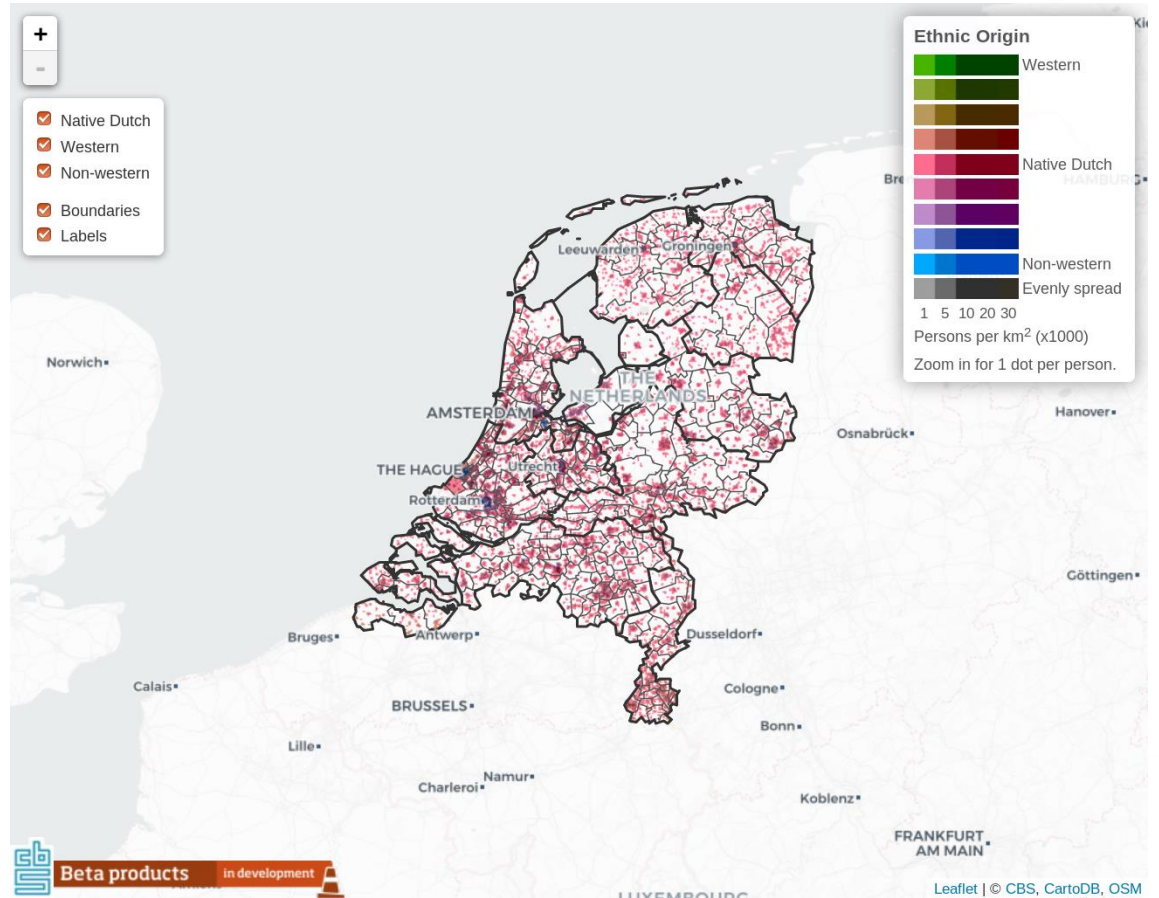


Published a CBDS beta product: https://research.cbs.nl/colordotmap/en

# **Application**

Migration background of
the Dutch population

Experimental version:
- Dots are placed in
  building areas
  (using the BAG register)
- "Artistic" legend

# **Application**

Migration background of the Dutch population

Experimental version:
- Dots are placed in building areas (using the BAG register)
- "Artistic" legend

# Application

Migration background of the Dutch population

Experimental version:
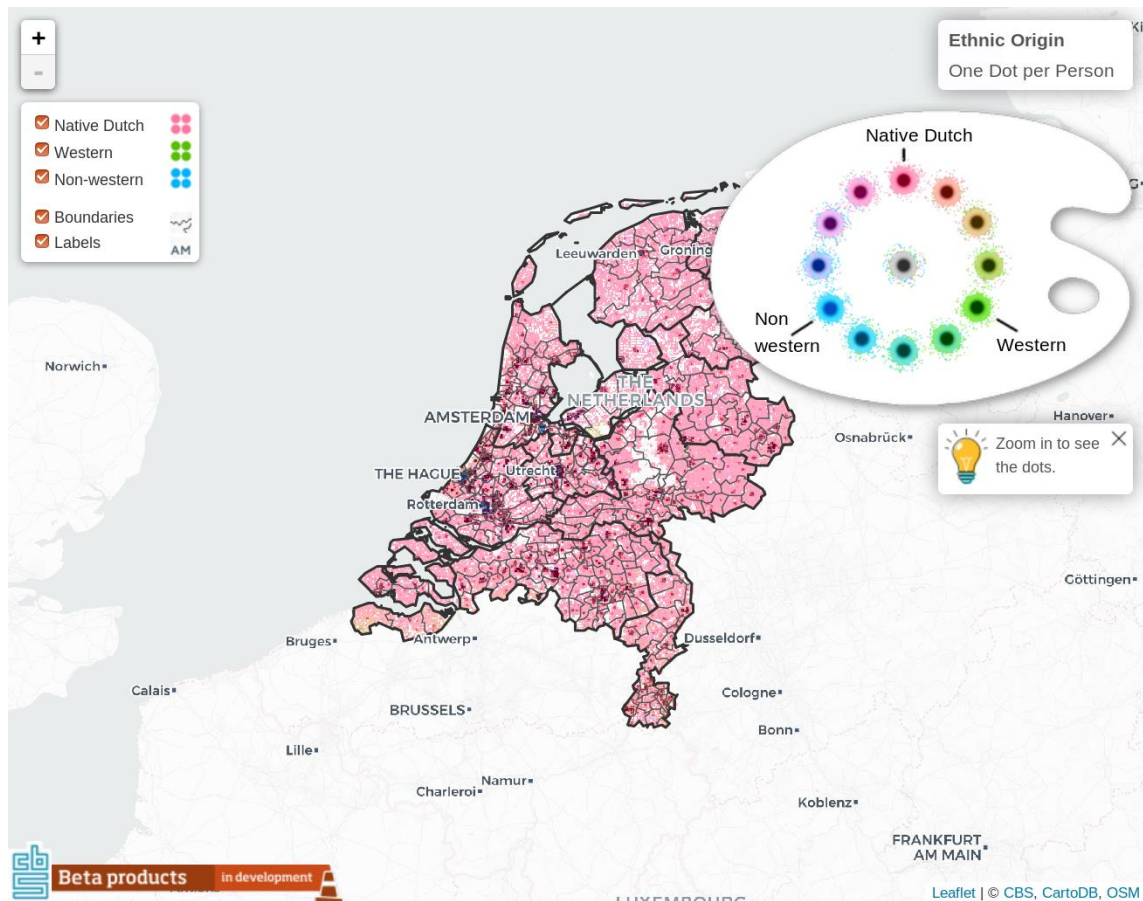- Dots are placed in building areas
  (using the BAG register)
- "Artistic" legend
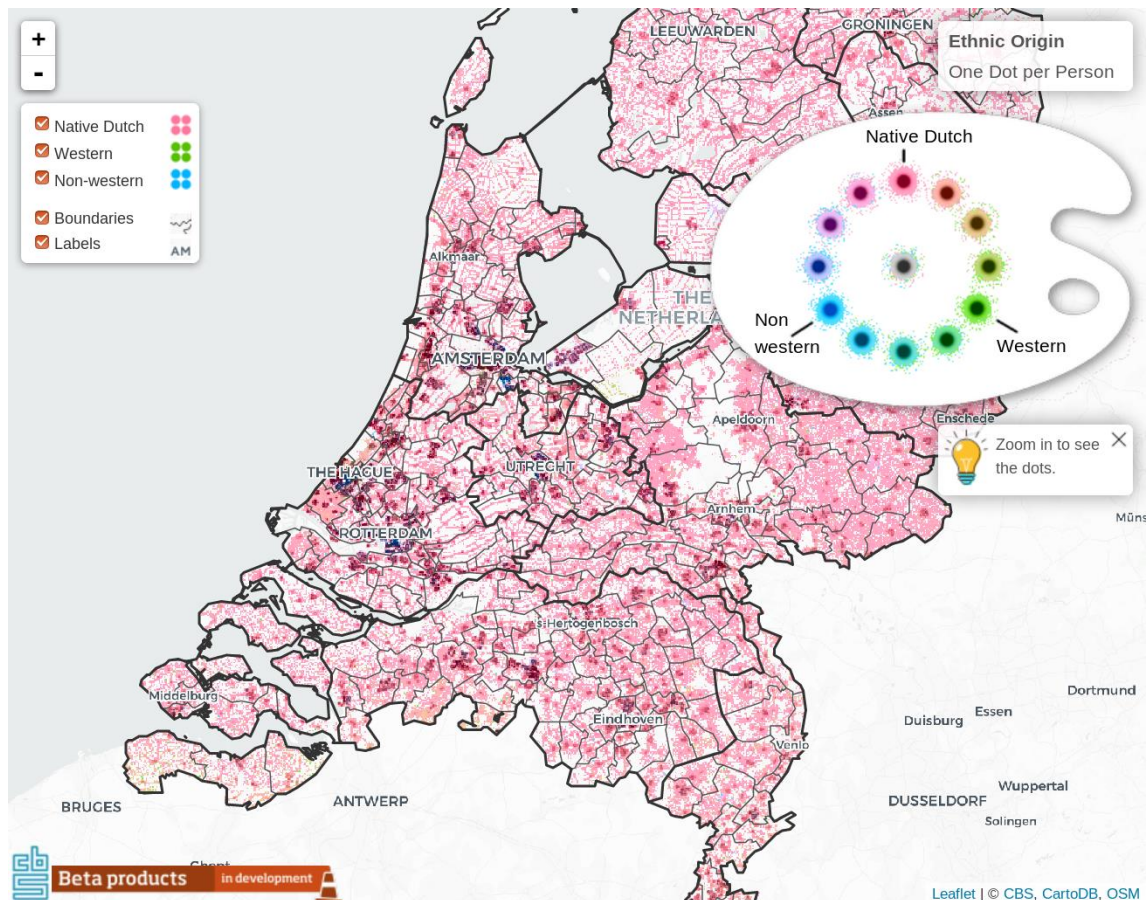
# Application

Migration background of the Dutch population

Experimental version:
- Dots are placed in building areas (using the BAG register)
- "Artistic" legend

# **Application**

Migration background of the Dutch population

Experimental version:
- Dots are placed in building areas (using the BAG register)
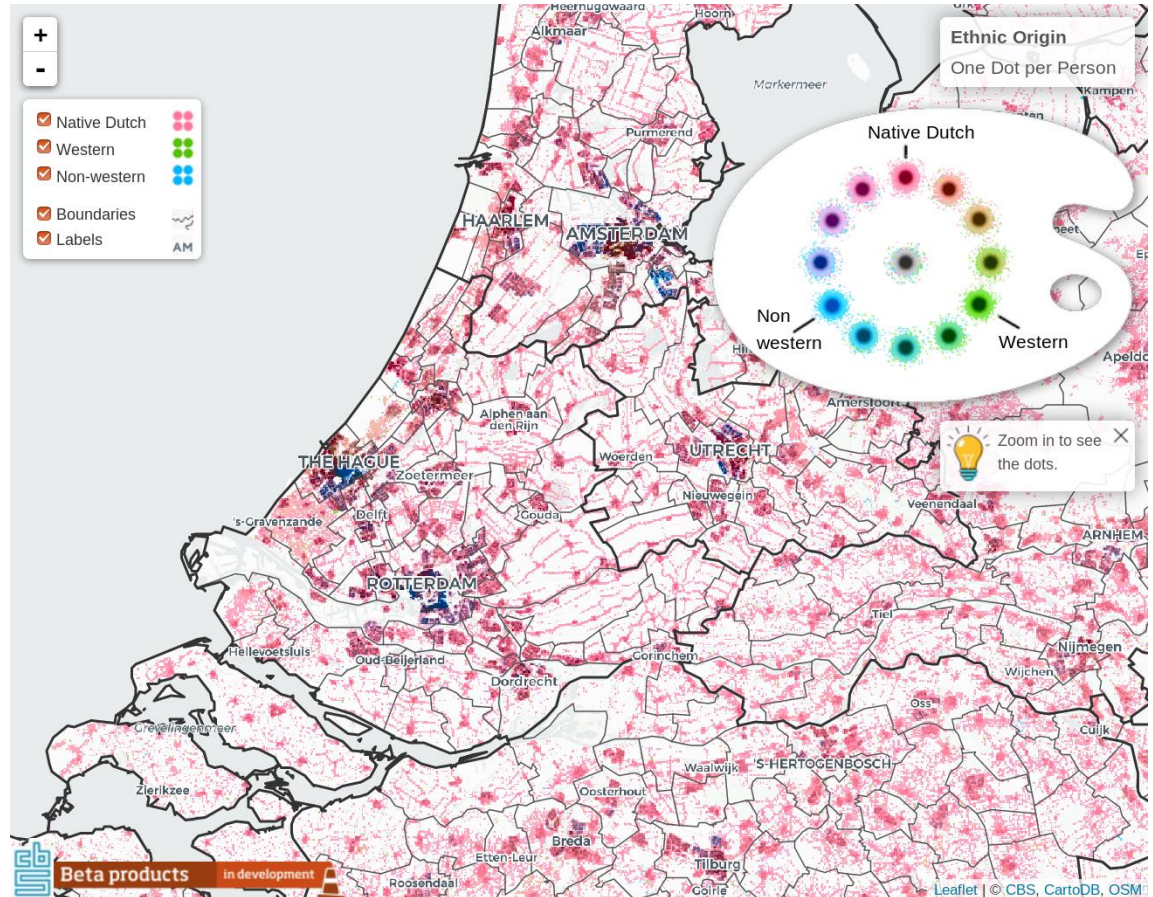- "Artistic" legend

# **Application**

Migration background of the Dutch population

Experimental version:
- Dots are placed in building areas (using the BAG register)
- "Artistic" legend

# **Application**

Migration background of
the Dutch population

Experimental version:
- Dots are placed in
  building areas
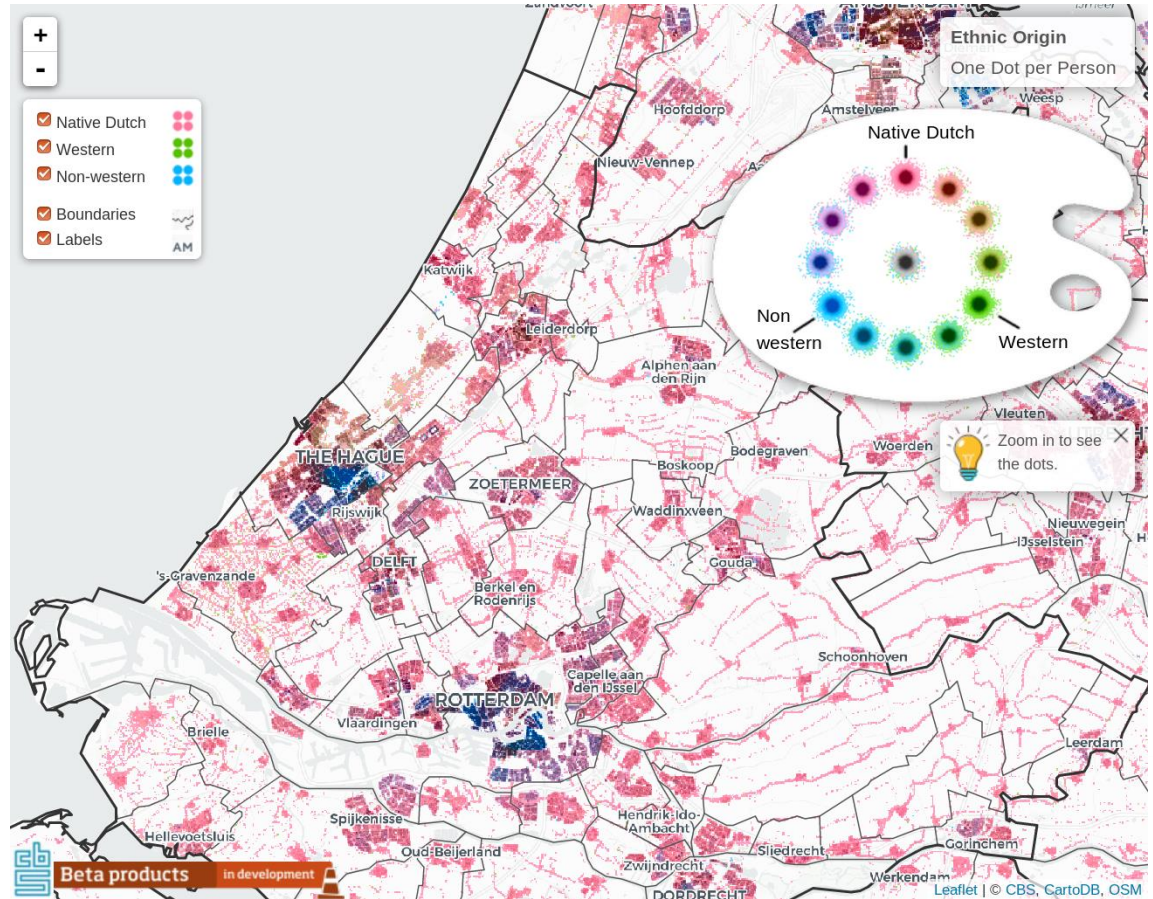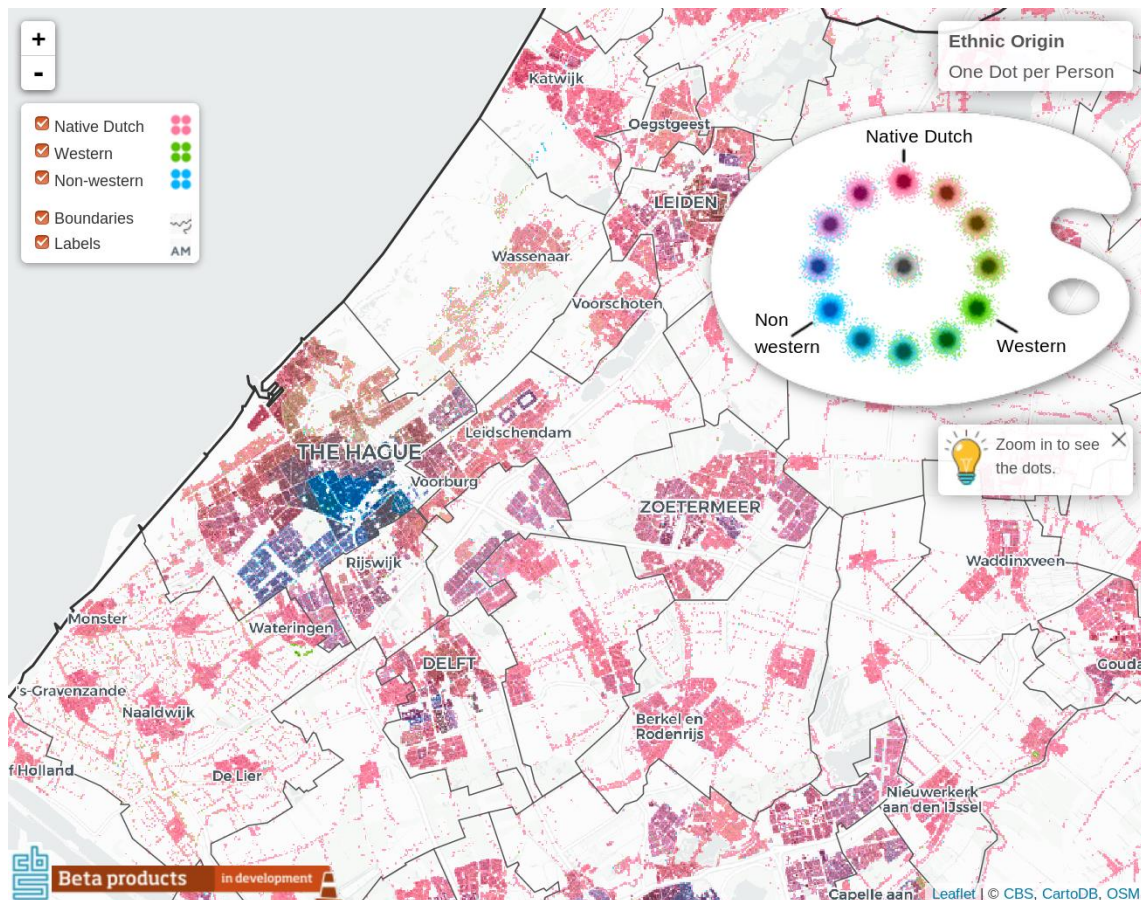  (using the BAG register)
- "Artistic" legend

# Application

Migration background of
the Dutch population

Experimental version:
- Dots are placed in
  building areas
  (using the BAG register)
- "Artistic" legend
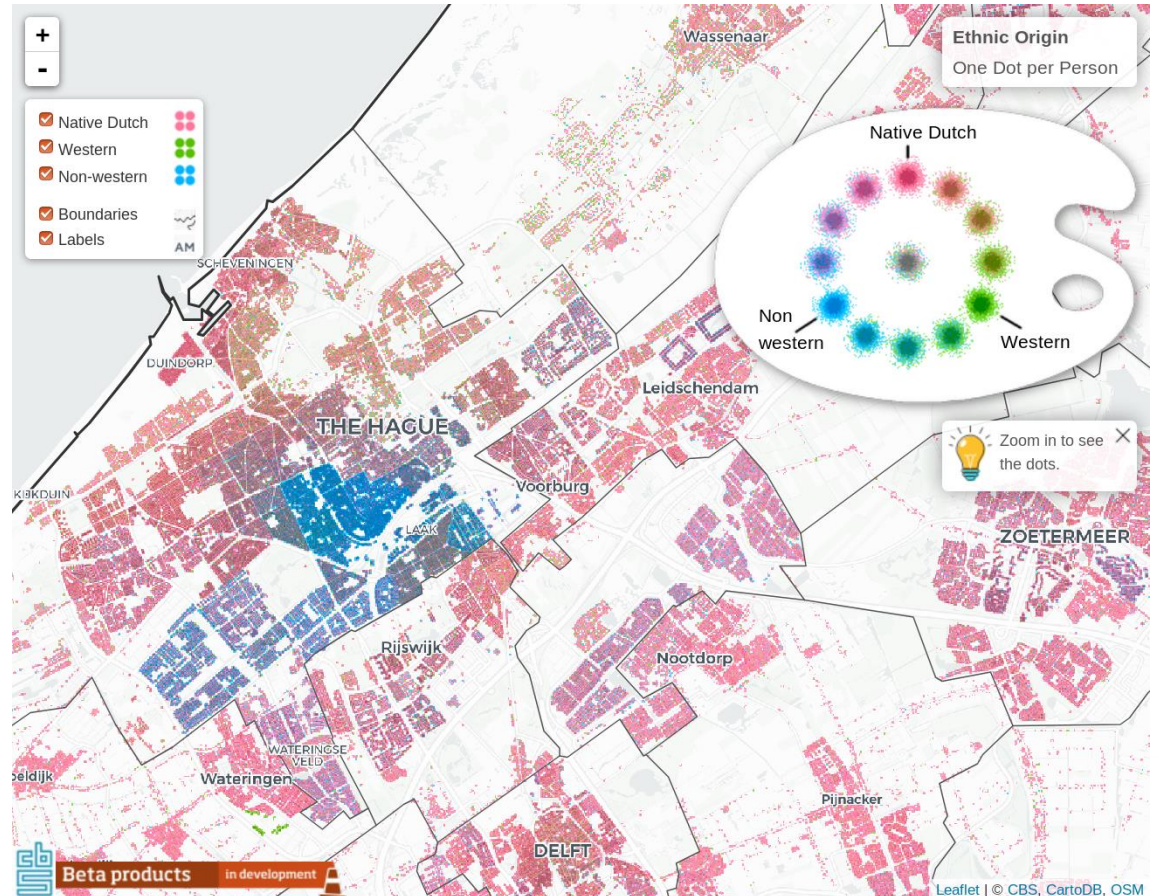
# Application

Migration background of the Dutch population

Experimental version:
- Dots are placed in building areas
  (using the BAG register)
- "Artistic" legend

# Application

Migration background of the Dutch population

Experimental version:
- Dots are placed in building areas
  (using the BAG register)
- "Artistic" legend

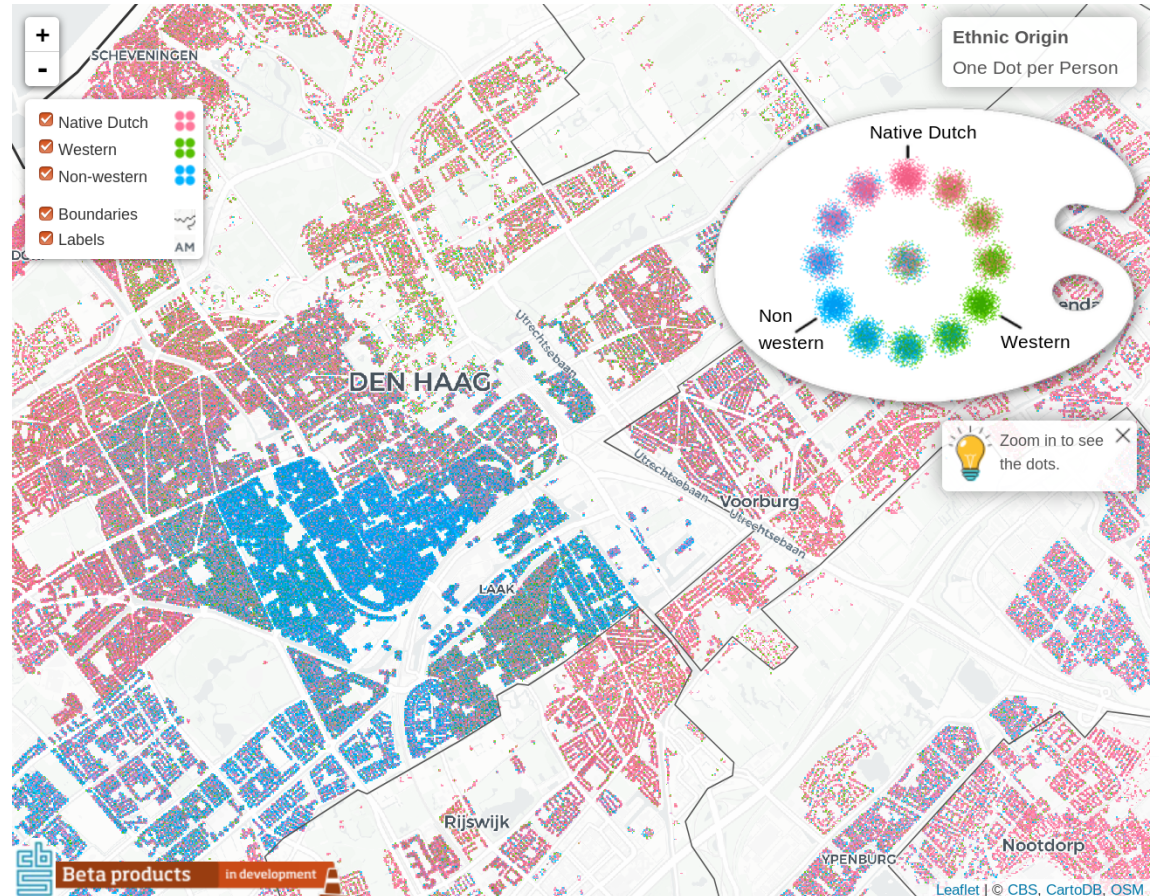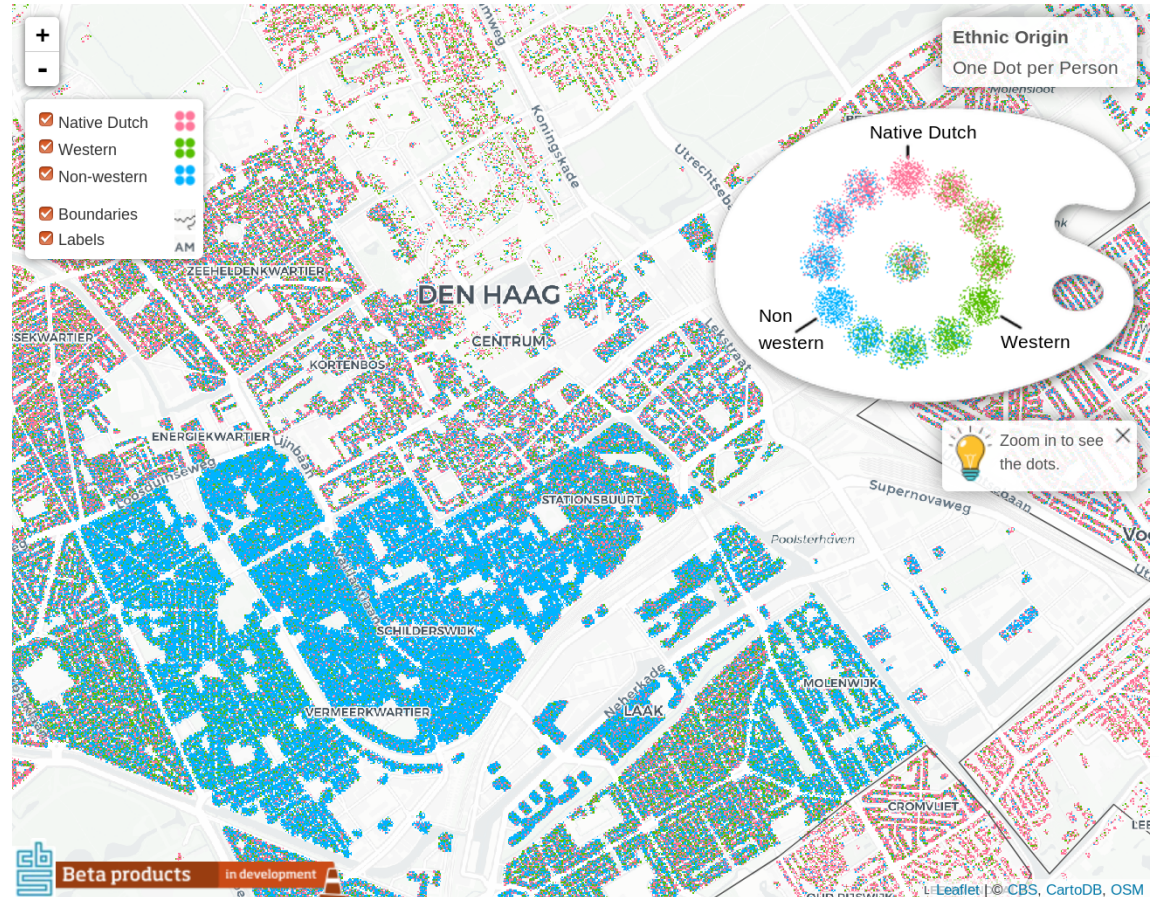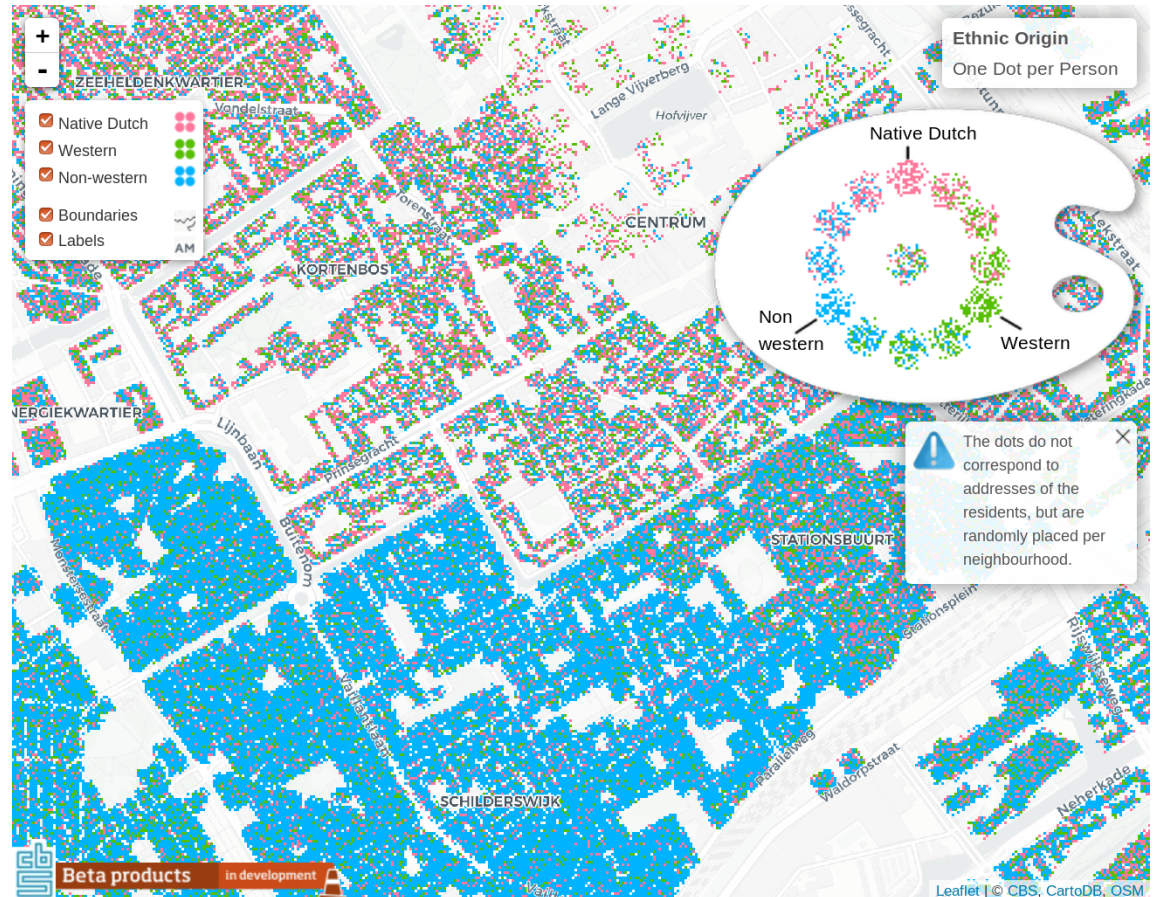# **Application**

Migration background of the Dutch population

Experimental version:
- Dots are placed in building areas (using the BAG register)
- "Artistic" legend

# User study

Comparison between original and experimental version with eye-tracking.



*Strange… Neighbourhoods appear pink from a distance, but from nearby, you clearly see the mix.*

# User study

Comparison between **original** and **experimental version** with **eye-tracking**.



*Pink is dominant, and therefore it's hard to distinguish between green and blue.*

# User study

Comparison between **original** and **experimental version** with **eye-tracking**.

# User study

**Conclusion:**

- Discrepancy between nearby and distant views, although users were able to read and interpret composition and density correctly.
- Legend was difficult to interpret (both versions).
- Most users thought that the dots where placed on actual addresses.

# How to deal with privacy?

Some ideas / guidelines:

- Areas should not be too detailed (global land use is better than detailed building areas)
- Draw neighbourhood borders
- Limit the zoom level (not to close)

# Application



- Simulated data on neighbourhood level for Amsterdam
- Each dot represents a household
- Dots are placed in residential areas (OpenStreetMap) per neighbourhood

**Welcome to ClairCity**

Citizen-led air pollution reduction in cities

WHERE IS
CLAIRCITY?



AMSTERDAM

BRISTOL

SOSNOWIEC

AVEIRO

LIGURIA

LJUBLJANA

http://www.claircity.eu/

34

# **Application**



ClairCity

- Simulated data on neighbourhood level for Amsterdam
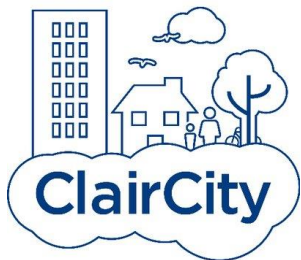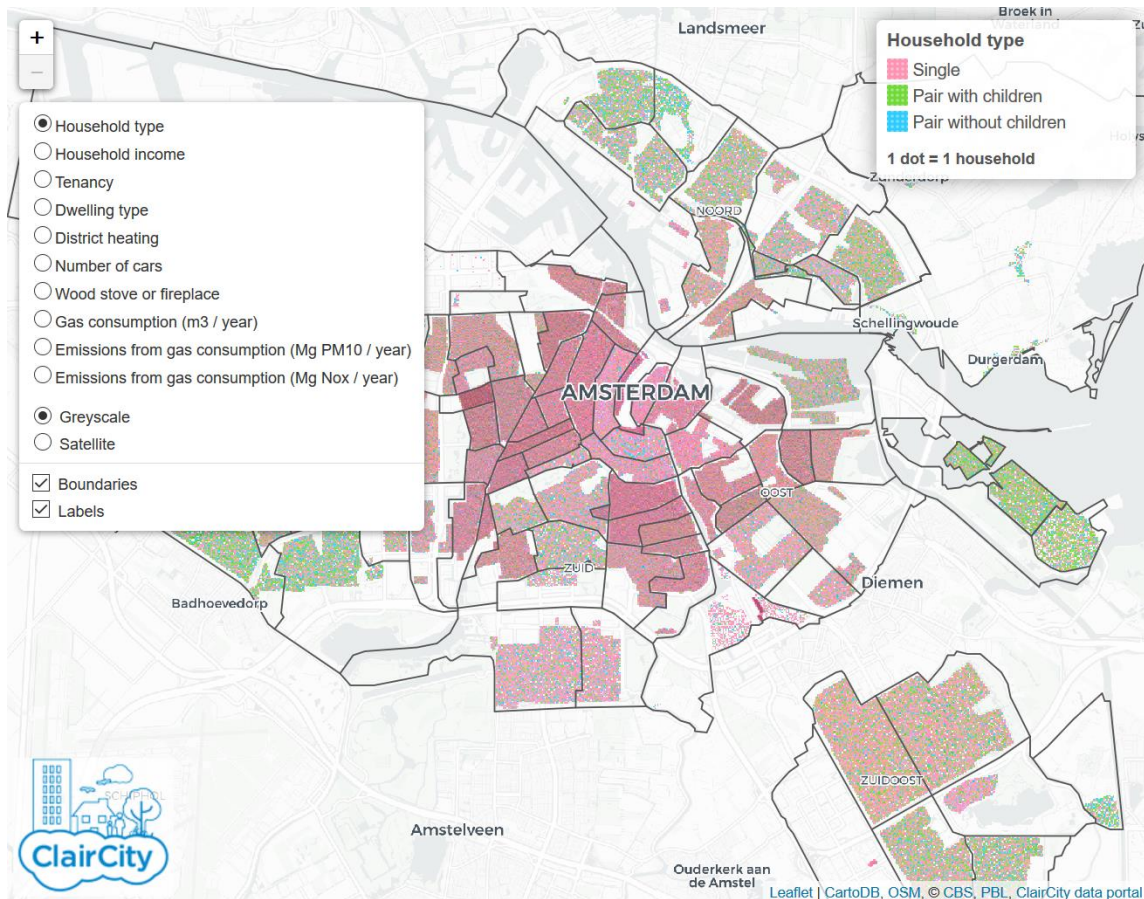- Each dot represents a household
- Dots are placed in residential areas (OpenStreetMap) per neighbourhood



Household type
- Single
- Pair with children
- Pair without children

1 dot = 1 household

○ Household type
○ Household income
○ Tenancy
○ Dwelling type
○ District heating
○ Number of cars
○ Wood stove or fireplace
○ Gas consumption (m3 / year)
○ Emissions from gas consumption (Mg PM10 / year)
○ Emissions from gas consumption (Mg Nox / year)

● Greyscale
○ Satellite

☑ Boundaries
☑ Labels

http://www.claircity.eu/
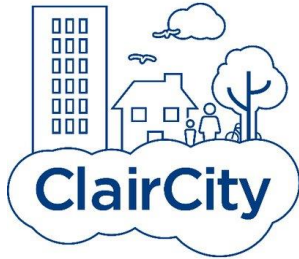
# Application


ClairCity

- Simulated data on neighbourhood level for Amsterdam
- Each dot represents a household
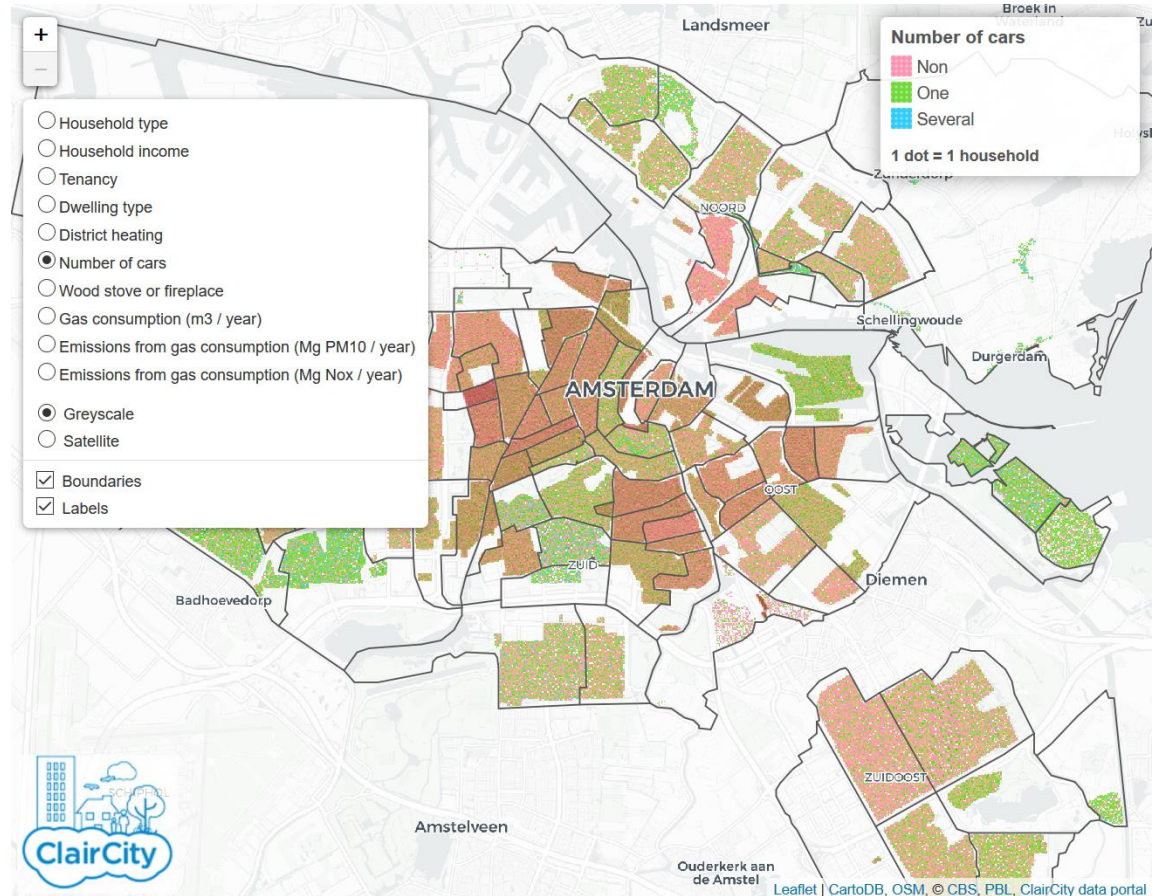- Dots are placed in residential areas (OpenStreetMap) per neighbourhood



- Household type
- Household income
- Tenancy
- Dwelling type
- District heating
- Number of cars (selected)
- Wood stove or fireplace
- Gas consumption (m3 / year)
- Emissions from gas consumption (Mg PM10 / year)
- Emissions from gas consumption (Mg Nox / year)

- Greyscale (selected)
- Satellite

- ☑ Boundaries
- ☑ Labels

**Number of cars**
- Non
- One
- Several

**1 dot = 1 household**

Leaflet | CartoDB, OSM, © CBS, PBL, ClairCity data portal

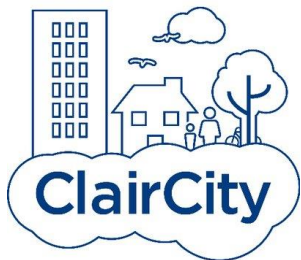http://www.claircity.eu/

36

# Application



ClairCity

- Simulated data on neighbourhood level for Amsterdam
- Each dot represents a household
- Dots are placed in residential areas (OpenStreetMap) per neighbourhood
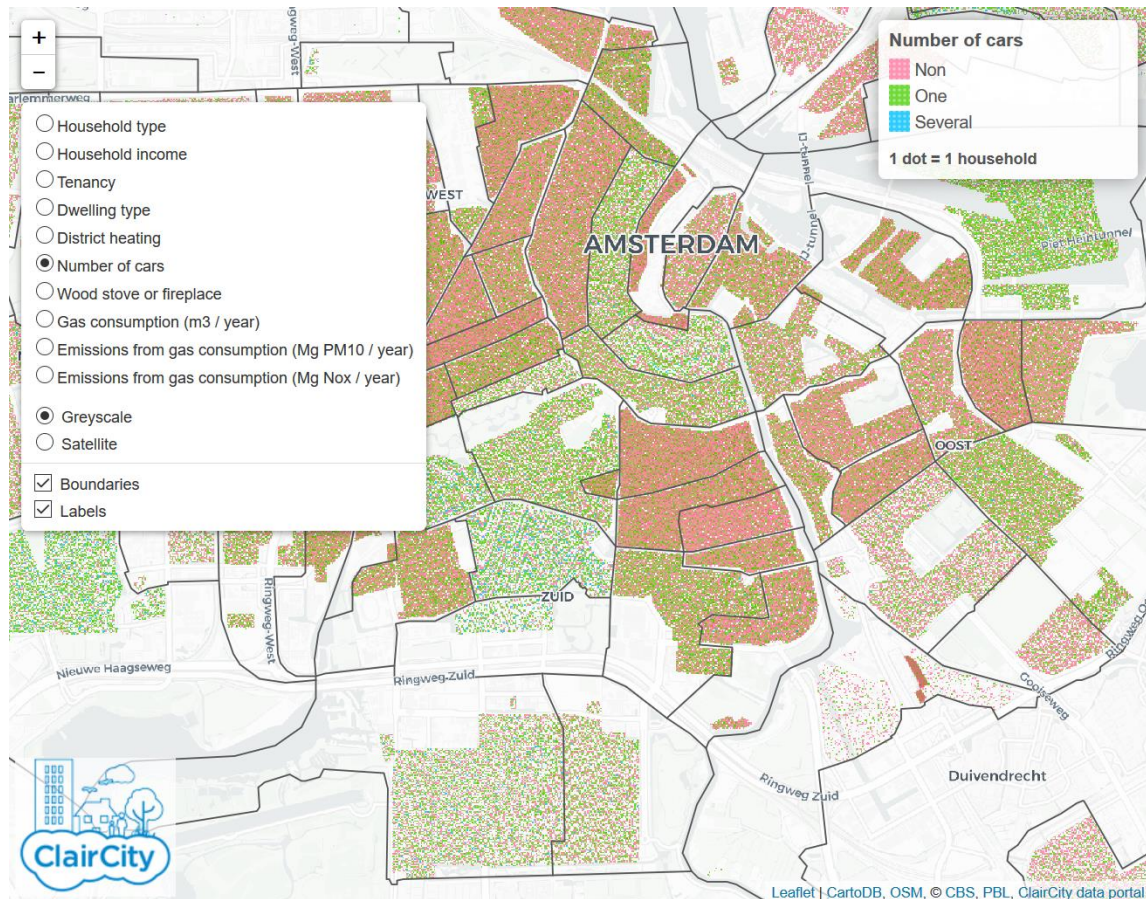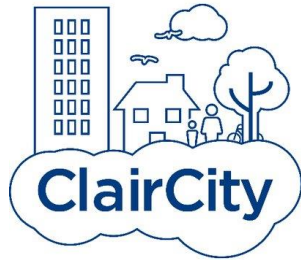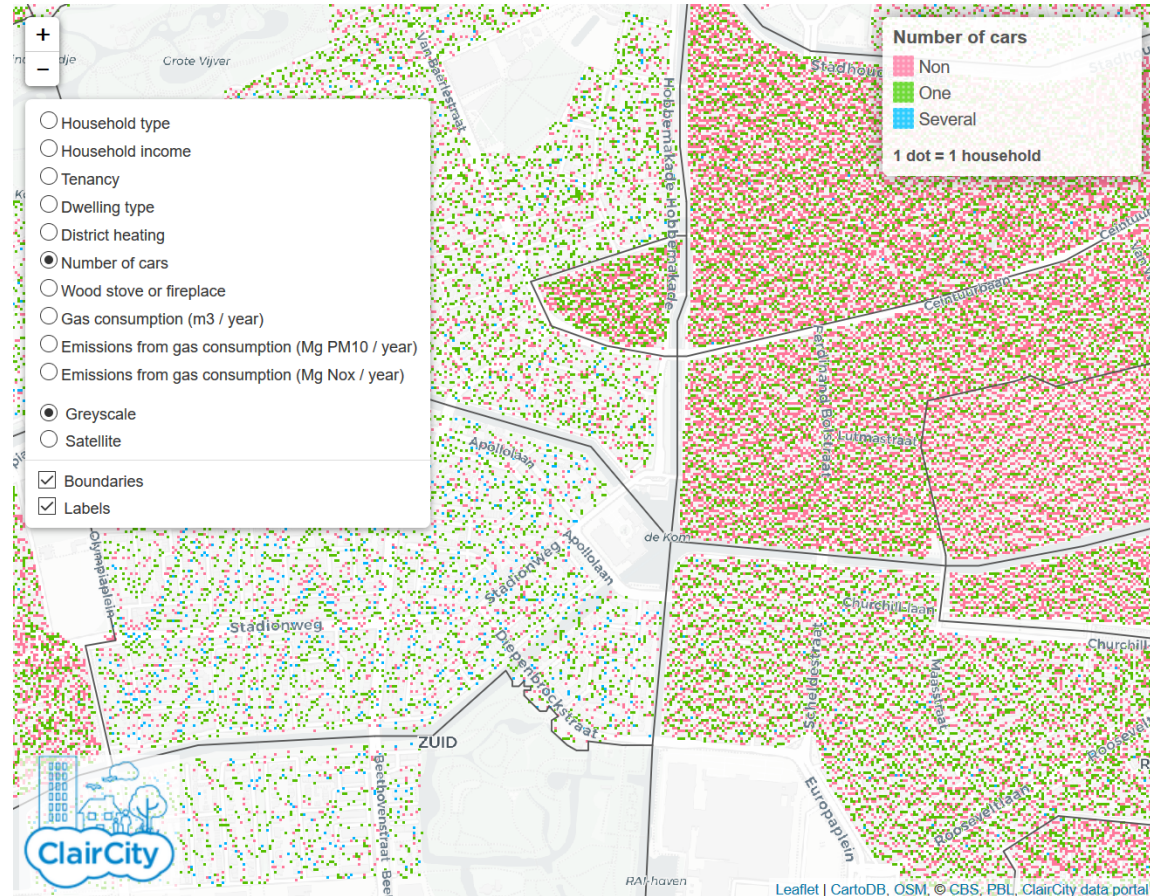


- Household type
- Household income
- Tenancy
- Dwelling type
- District heating
- ⦿ Number of cars
- Wood stove or fireplace
- Gas consumption (m3 / year)
- Emissions from gas consumption (Mg PM10 / year)
- Emissions from gas consumption (Mg Nox / year)

- ⦿ Greyscale
- Satellite

- ☑ Boundaries
- ☑ Labels

**Number of cars**
- Non
- One
- Several

1 dot = 1 household

Leaflet | CartoDB, OSM, © CBS, PBL, ClairCity data portal
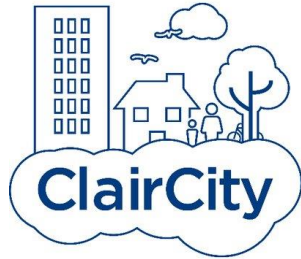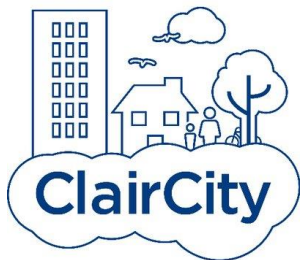
http://www.claircity.eu/

# Application



- Simulated data on neighbourhood level for Amsterdam
- Each dot represents a household
- Dots are placed in residential areas (OpenStreetMap) per neighbourhood

http://www.claircity.eu/

# Application



ClairCity

- Simulated data on neighbourhood level for Amsterdam
- Each dot represents a household
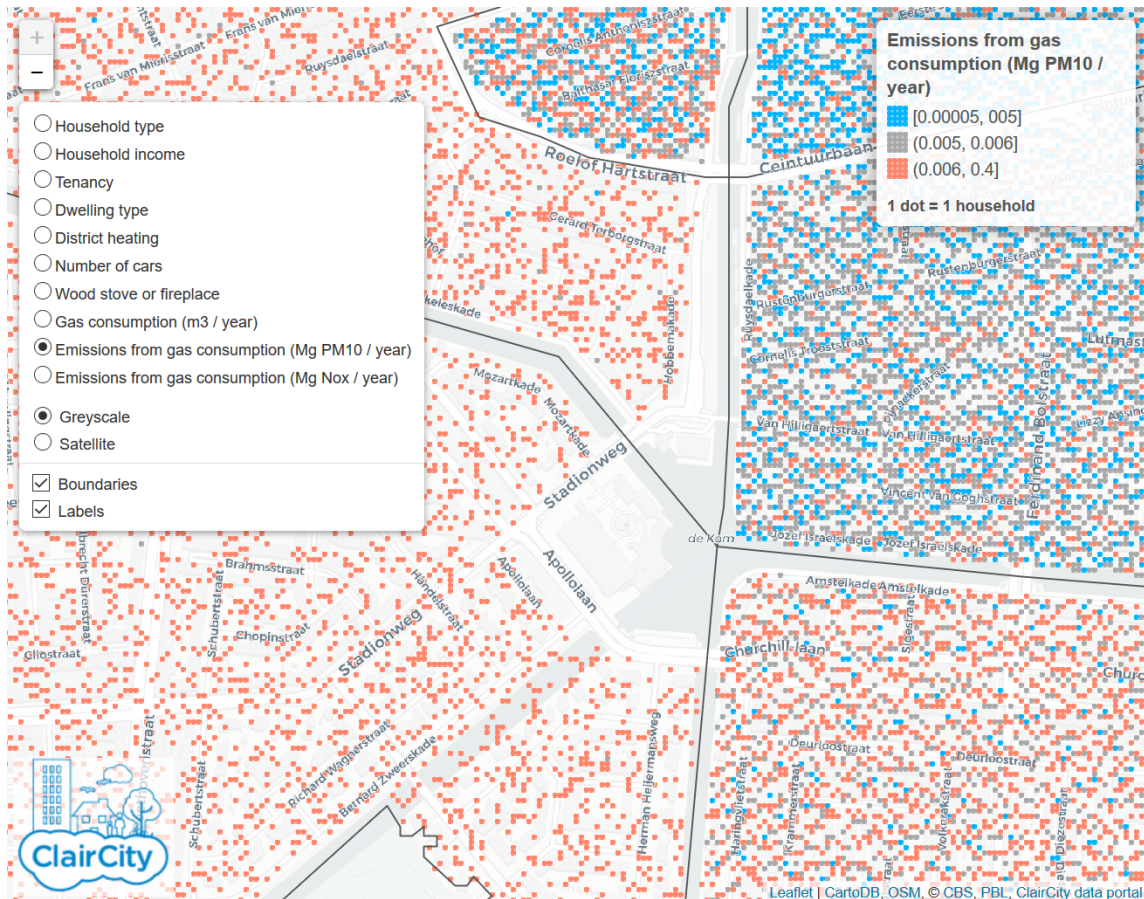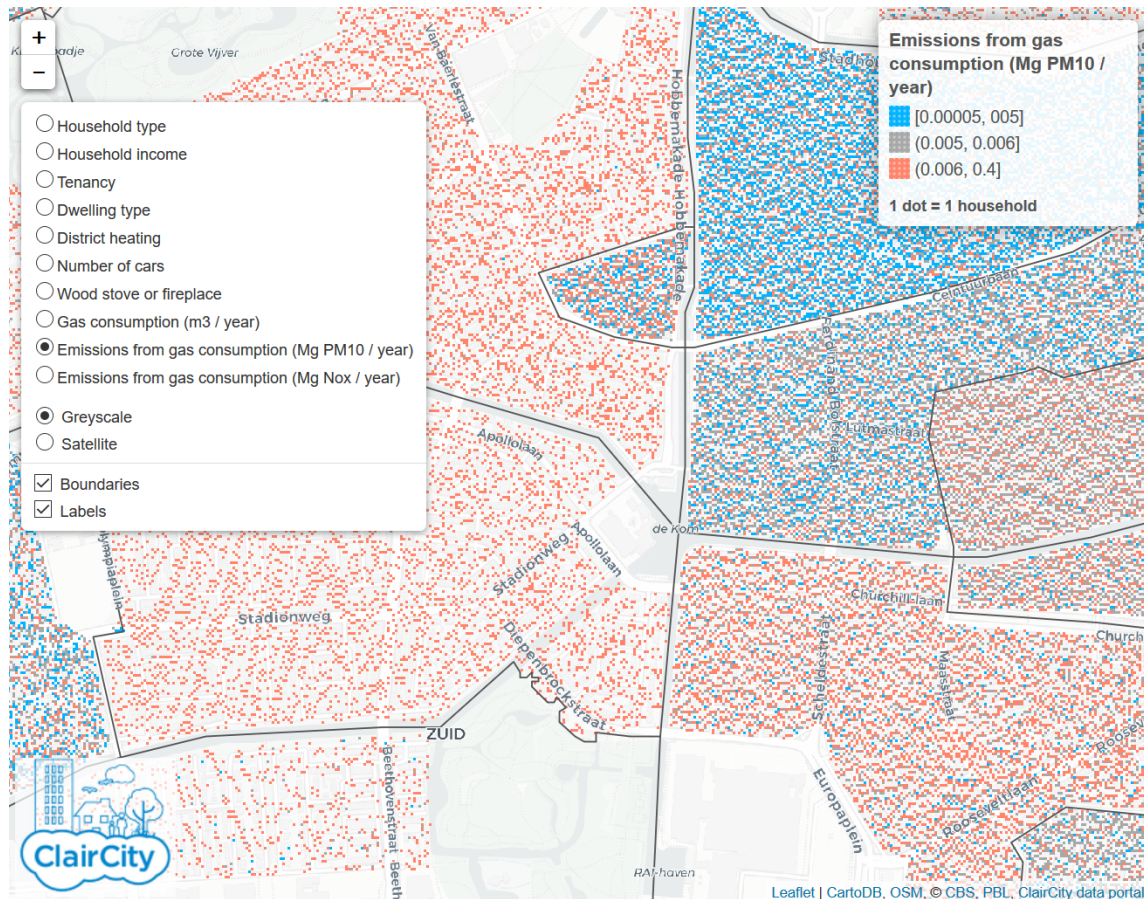- Dots are placed in residential areas (OpenStreetMap) per neighbourhood



**Number of cars**
- Non
- One
- Several

1 dot = 1 household

Household type
Household income
Tenancy
Dwelling type
District heating
● Number of cars
Wood stove or fireplace
Gas consumption (m3 / year)
Emissions from gas consumption (Mg PM10 / year)
Emissions from gas consumption (Mg Nox / year)

● Greyscale
Satellite

☑ Boundaries
☑ Labels

Leaflet | CartoDB, OSM, © CBS, PBL, ClairCity data portal
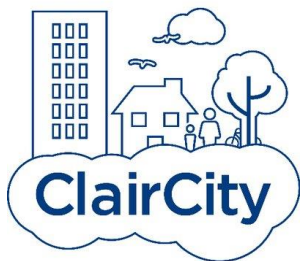
http://www.claircity.eu/

# **Application**



- Simulated data on neighbourhood level for Amsterdam
- Each dot represents a household
- Dots are placed in residential areas (OpenStreetMap) per neighbourhood
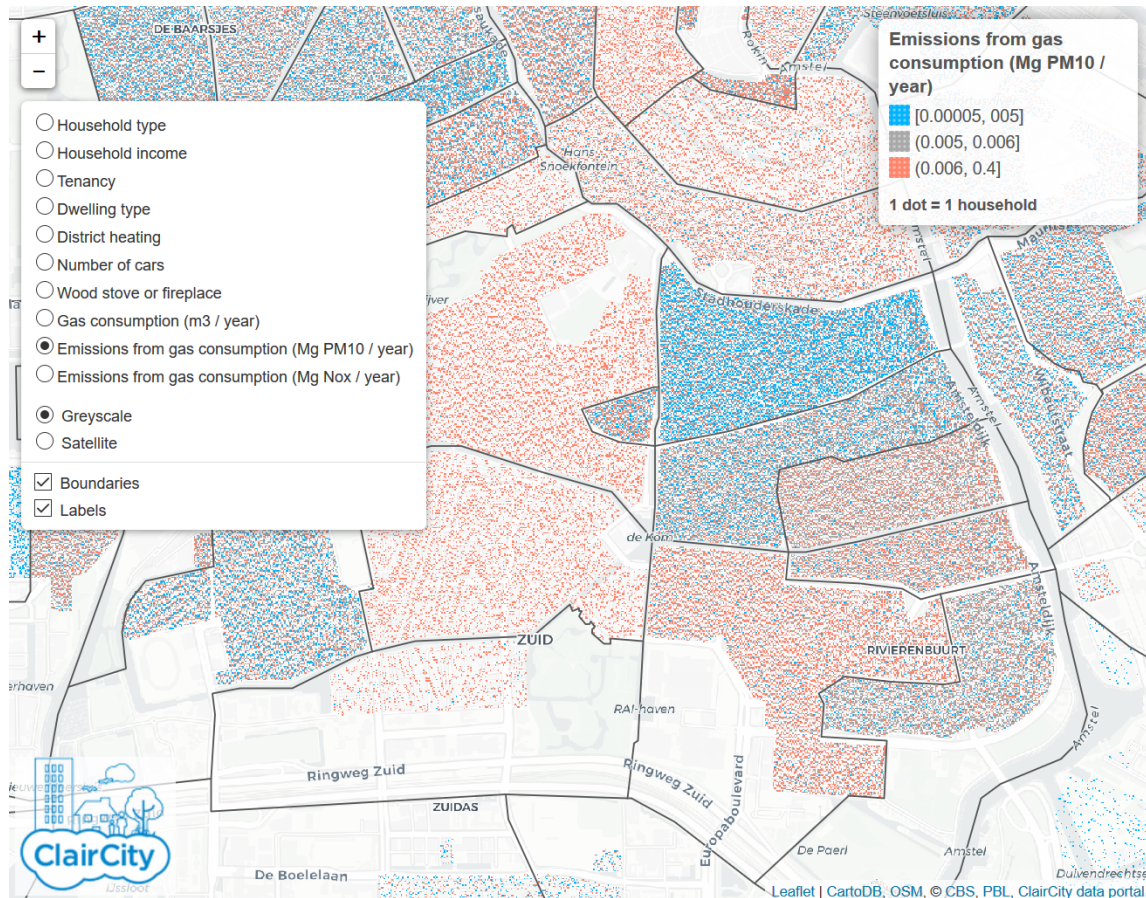


http://www.claircity.eu/
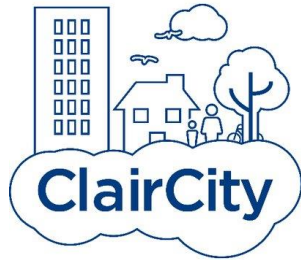
40

# **Application**



- Simulated data on neighbourhood level for Amsterdam
- Each dot represents a household
- Dots are placed in residential areas (OpenStreetMap) per neighbourhood

http://www.claircity.eu/

# **Application**



- Simulated data on neighbourhood level for Amsterdam
- Each dot represents a household
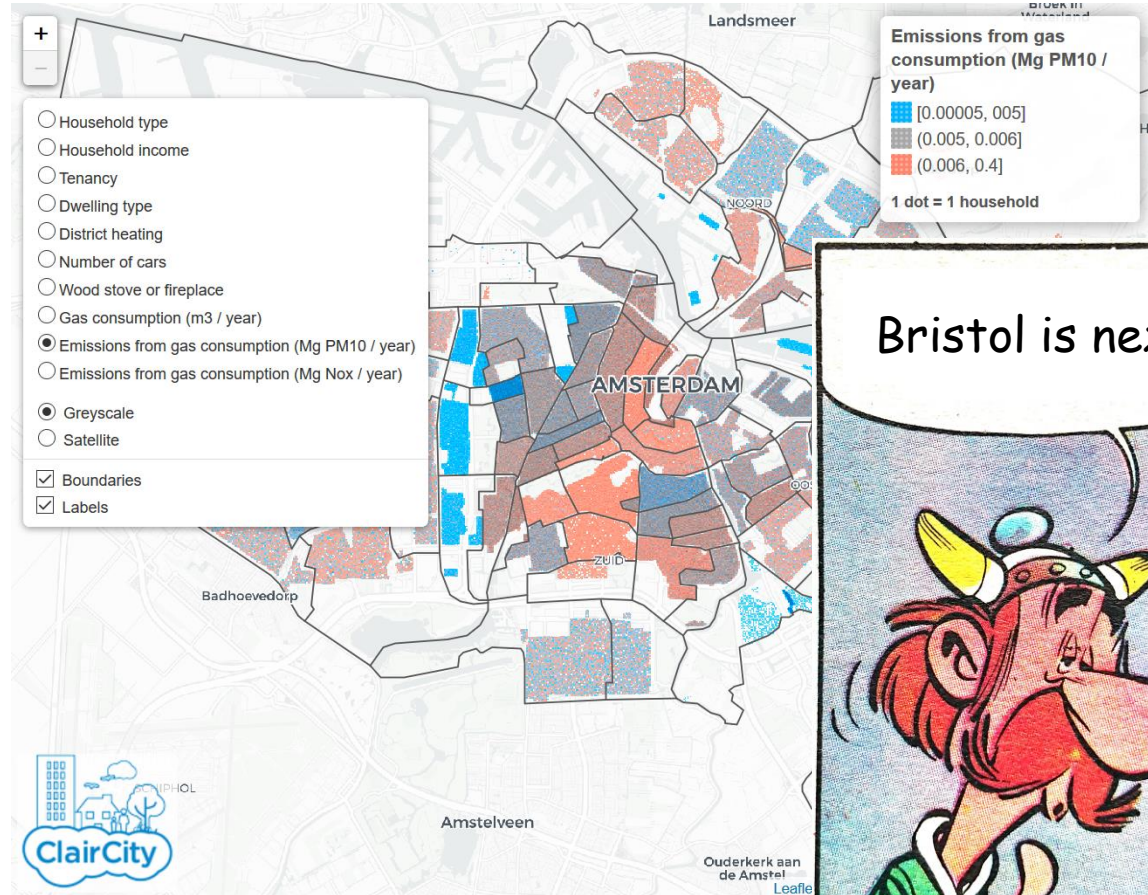- Dots are placed in residential areas (OpenStreetMap) per neighbourhood

http://www.claircity.eu/
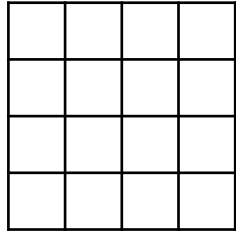
# **Application**



- Simulated data on neighbourhood level for Amsterdam
- Each dot represents a household
- Dots are placed in residential areas (OpenStreetMap) per neighbourhood
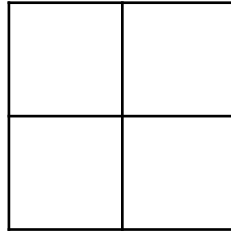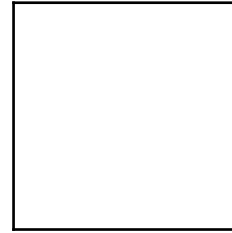


○ Household type
○ Household income
○ Tenancy
○ Dwelling type
○ District heating
○ Number of cars
○ Wood stove or fireplace
○ Gas consumption (m3 / year)
● Emissions from gas consumption (Mg PM10 / year)
○ Emissions from gas consumption (Mg Nox / year)

● Greyscale
○ Satellite

☑ Boundaries
☑ Labels

**Emissions from gas consumption (Mg PM10 / year)**
▦ [0.00005, 005]
▦ (0.005, 0.006]
▦ (0.006, 0.4]

1 dot = 1 household

Leaflet | CartoDB, OSM, © CBS, PBL, ClairCity data portal

http://www.claircity.eu/

43

# Application


ClairCity

- Simulated data on neighbourhood level for Amsterdam
- Each dot represents a household
- Dots are placed in residential areas (OpenStreetMap) per neighbourhood



Emissions from gas consumption (Mg PM10 / year)
- [0.00005, 005]
- (0.005, 0.006)
- (0.006, 0.4]

1 dot = 1 household

○ Household type
○ Household income
○ Tenancy
○ Dwelling type
○ District heating
○ Number of cars
○ Wood stove or fireplace
○ Gas consumption (m3 / year)
● Emissions from gas consumption (Mg PM10 / year)
○ Emissions from gas consumption (Mg Nox / year)

● Greyscale
○ Satellite

☑ Boundaries
☑ Labels



Bristol is next!

http://www.claircity.eu/

# Super Dots

*k* by *k* grid cells in **original matrix** = 1 grid cell in **aggregated matrix**



original          *k* = 2          *k* = 4

Example:



original dot map          aggregated dot map (*k* = 2)

# What is a good aggregation?



- **Class Balance** Total number of super dots per class should represent the total number of small dots per class

- **Representation** How well do the super dots represent the small dots? Each small dot is represented at most once, and each super dot can represent at most $k^2$ small dots.

- **Presence** How well are the small dots represented by the super dots? For each small dot, the distance to the nearest super dot is measured.
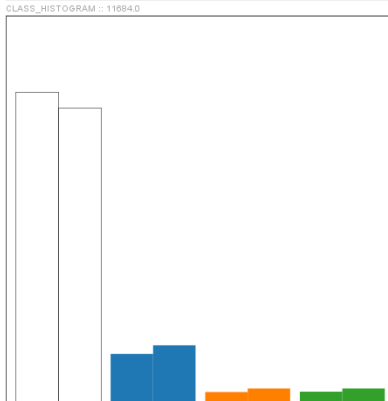
# Aggregation analyses tool



Original dot map

Aggregated dot map

Overlay

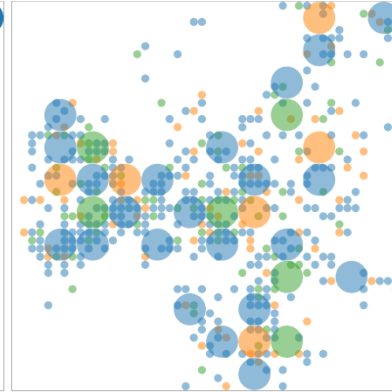Class balance

Representation

Presence

# Aggregation analyses tool



Original dot map

Aggregated dot map

Overlay

Class balance

Representation

Presence

# Aggregation analyses tool



Original dot map

Aggregated dot map

Overlay

Class balance

Representation

Presence

# Algorithms (sketches)

**Greedy Class Balance Algorithm**
1. Start with an empty map.
2. Pick the class with the largest imbalance and place a super dot of this class on the spot with the best representation.
3. Repeat step 2 until all super dots are placed.

**Kernel Density Sampling Algorithm**
1. For each class, estimate 2D kernel density.
2. Place super dots where total density is above a certain threshold.
3. Per super dot, sample its class using the density values as probabilities.

# Application



**Distance to school**

- Dots represent children who go to primary schools

- Colour indicates distance to their primary school (not necessarily the nearest one)

- Used data: education registers

- Draft version (not published yet)

- Dots aggregated using the Kernel Density Sampling Algorithm (only one aggregation)

# **Application**



**Distance to school**

- Dots represent children who go to primary schools

- Colour indicates distance to their primary school (not necessarily the nearest one)

- Used data: education registers

- Draft version (not published yet)

- Dots aggregated using the Kernel Density Sampling Algorithm (only one aggregation)
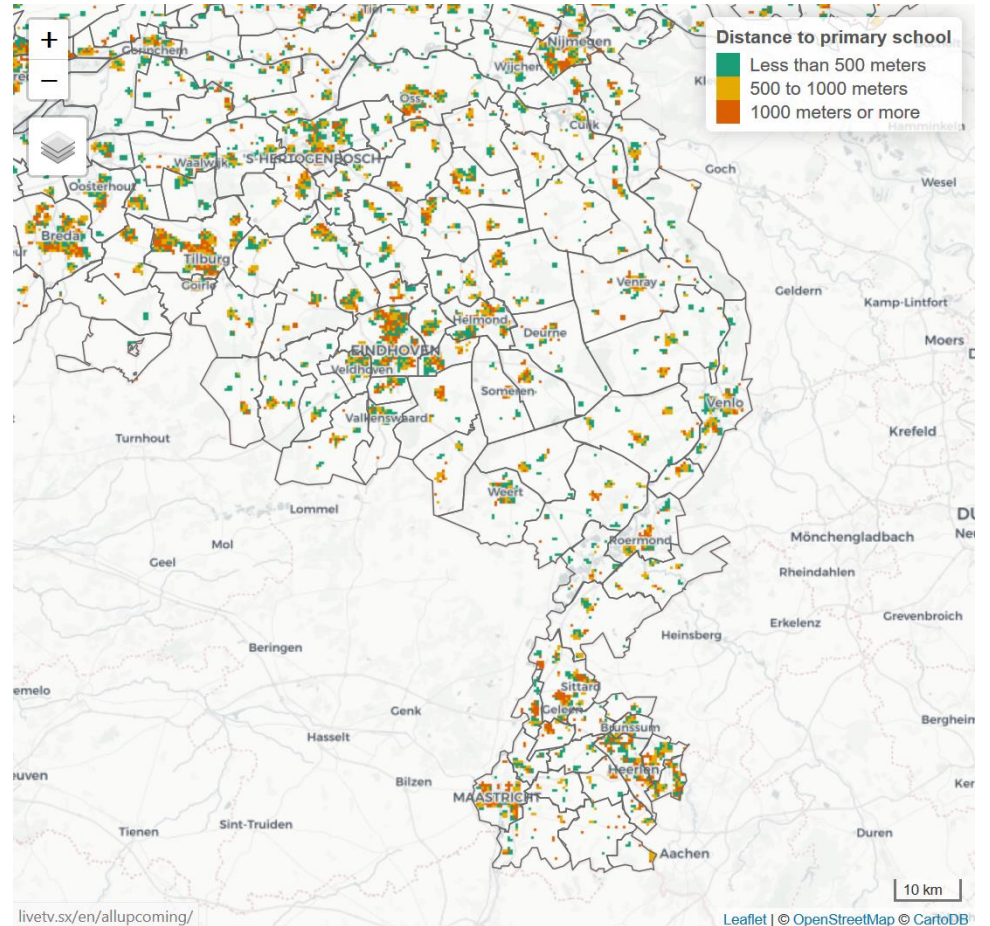
# Application

**Distance to school**

- Dots represent children who go to primary schools

- Colour indicates distance to their primary school (not necessarily the nearest one)

- Used data: education registers

- Draft version (not published yet)

- Dots aggregated using the Kernel Density Sampling Algorithm (only one aggregation)
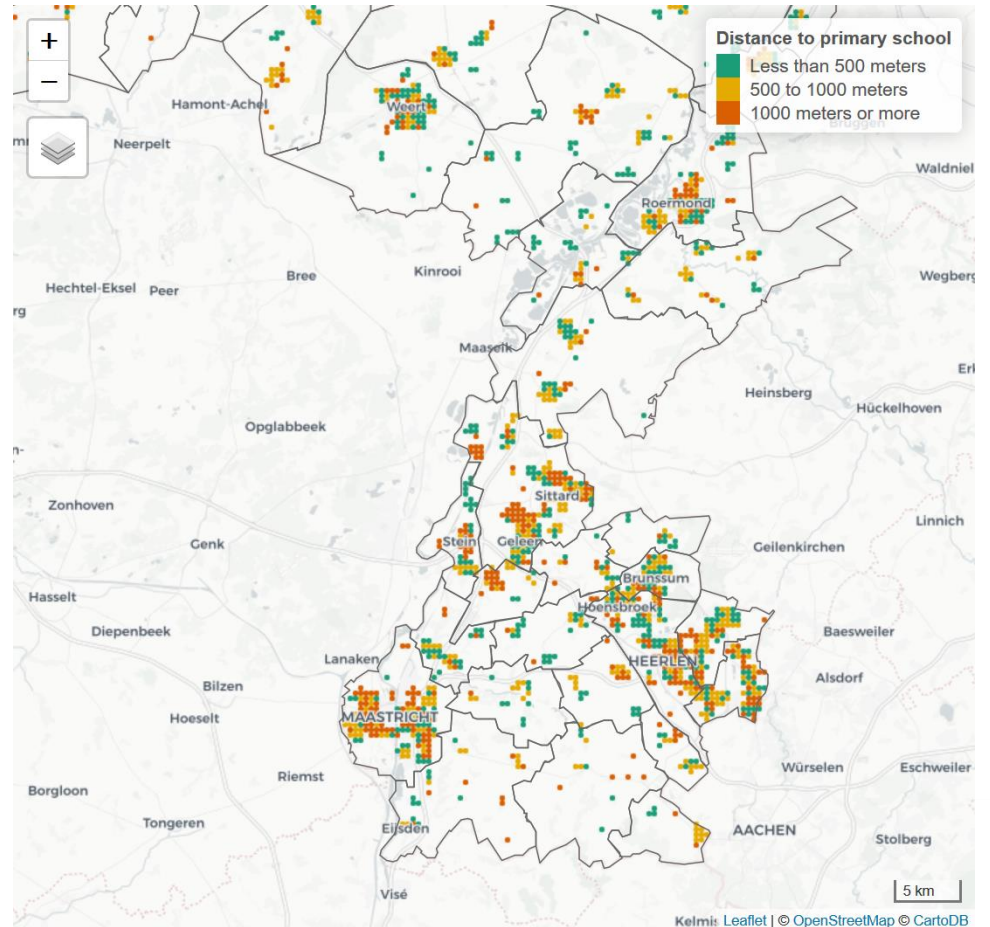


53

# Application



**Distance to school**

- Dots represent children who go to primary schools

- Colour indicates distance to their primary school (not necessarily the nearest one)

- Used data: education registers

- Draft version (not published yet)

- Dots aggregated using the Kernel Density Sampling Algorithm (only one aggregation) 54
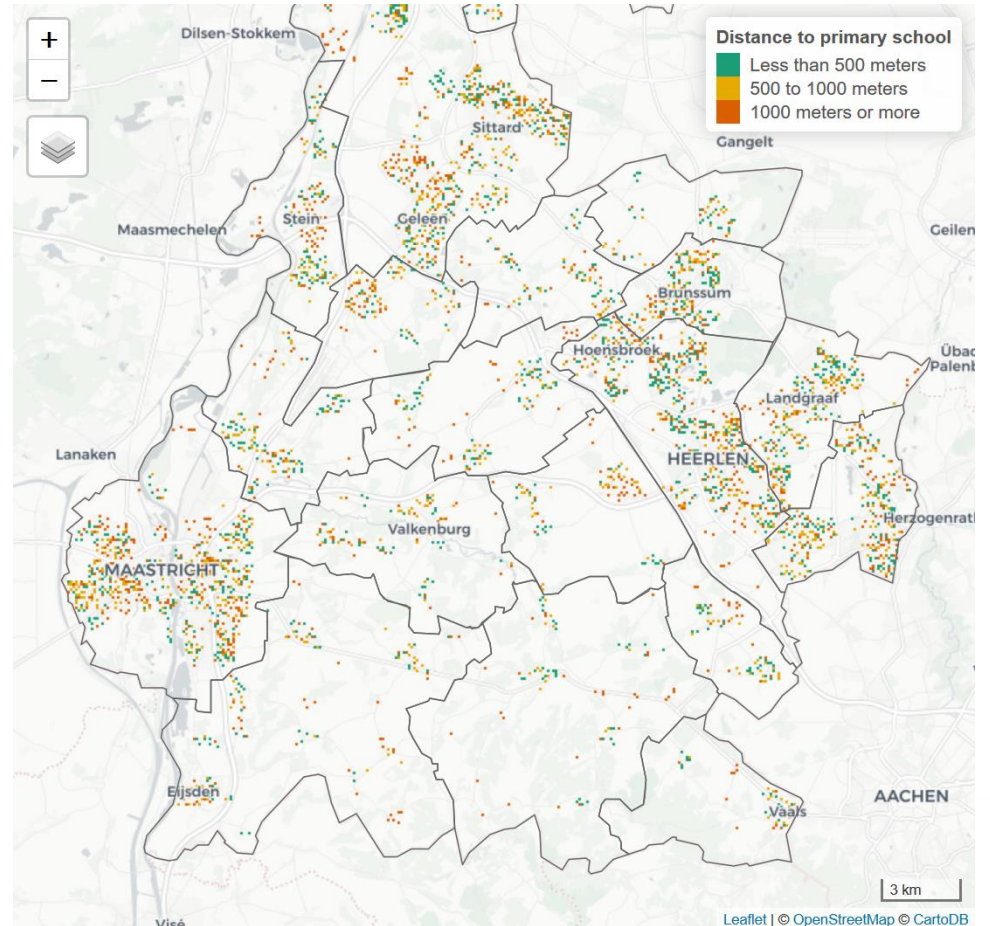
# Application



**Distance to school**

- Dots represent children who go to primary schools

- Colour indicates distance to their primary school (not necessarily the nearest one)

- Used data: education registers

- Draft version (not published yet)

- Dots aggregated using the Kernel Density Sampling Algorithm (only one aggregation)

# Application
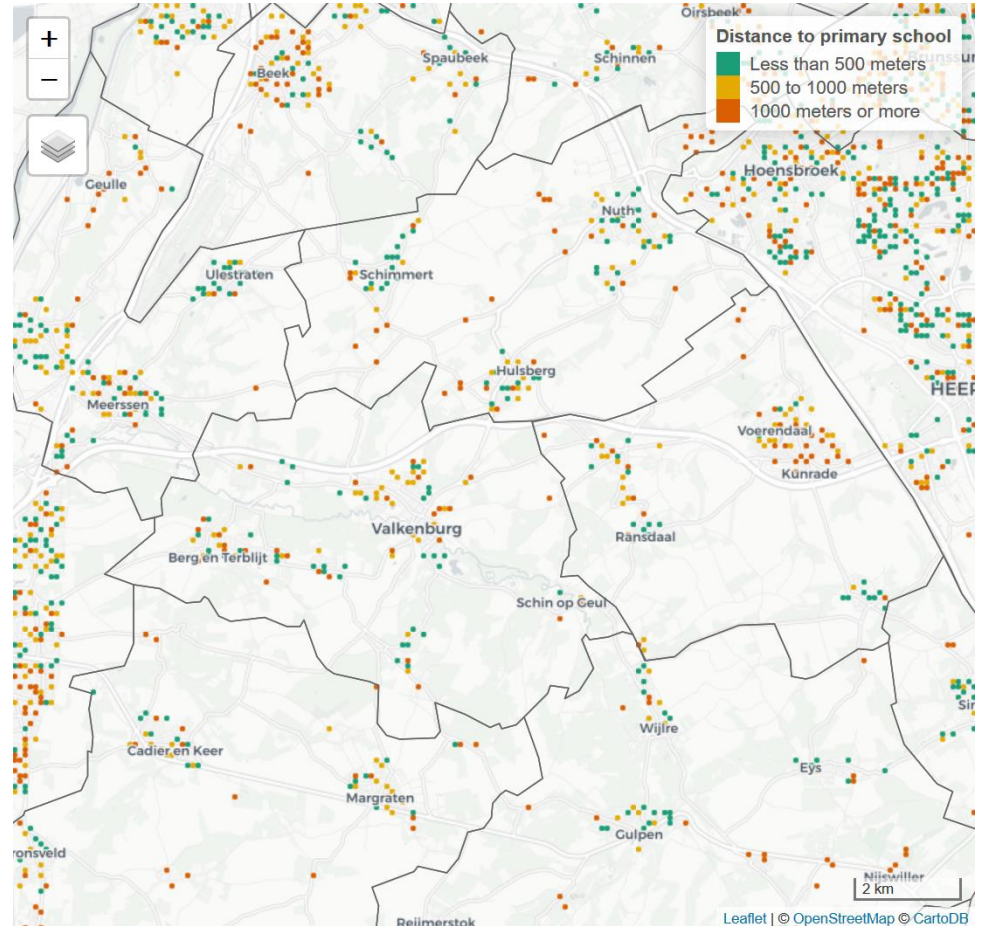


**Distance to school**

- Dots represent children who go to primary schools

- Colour indicates distance to their primary school (not necessarily the nearest one)

- Used data: education registers

- Draft version (not published yet)

- Dots aggregated using the Kernel Density Sampling Algorithm (only one aggregation)  56

# Application
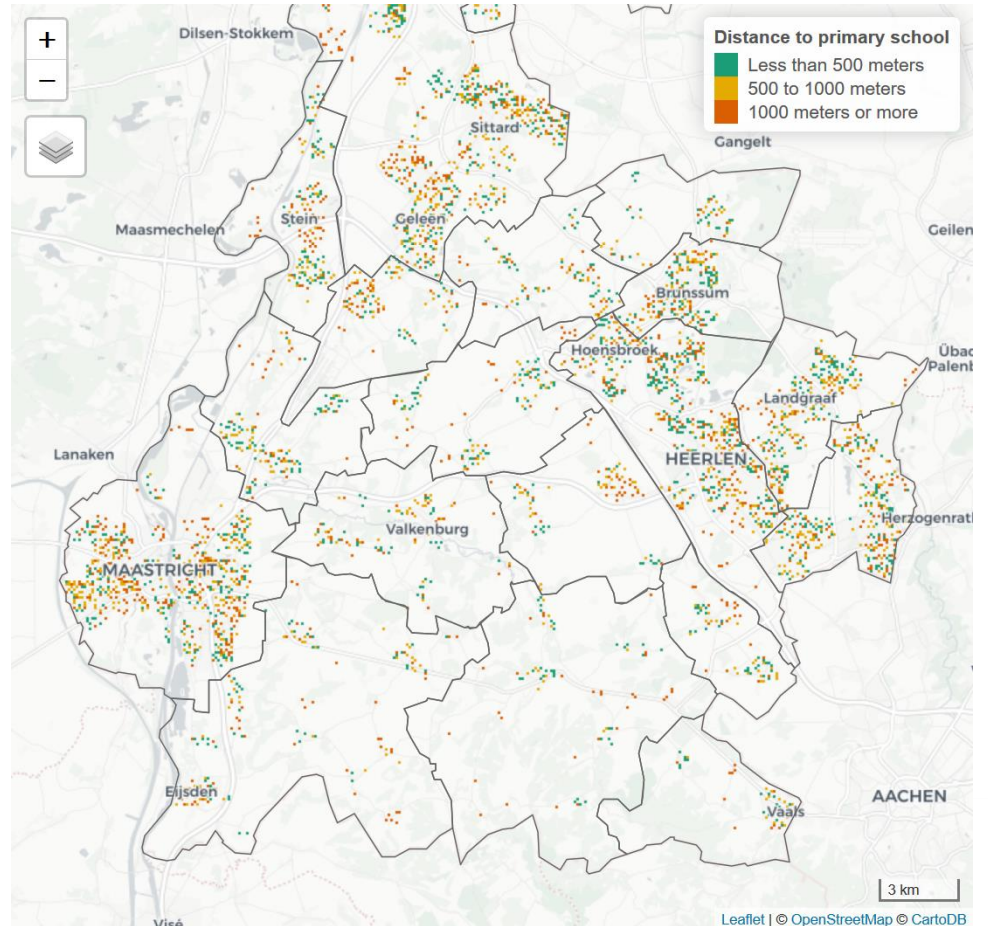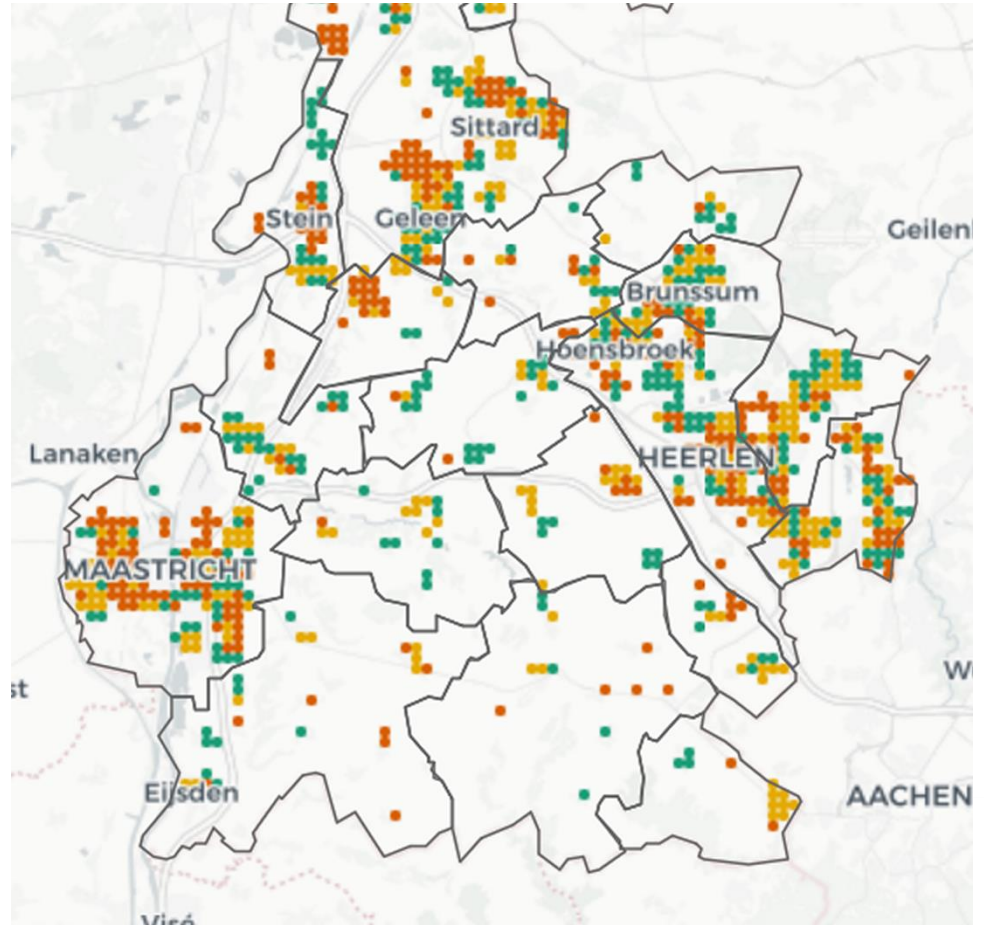
**Distance to school**

- Dots represent children who go to primary schools
- Colour indicates distance to their primary school (not necessarily the nearest one)
- Used data: education registers
- Draft version (not published yet)
- Dots aggregated using the Kernel Density Sampling Algorithm (only one aggregation)

# Comparison

**Blended colours**

+ Sense of immensity of the data

– Dots hard to distinguish and categorize

– Difficult to create simple legend

– Tricky to pick suitable colours (visual perception is complex)

**Super dots**

+ Simple and clear representation

+ Keeps the overall distribution and composition

– Loss of local detail

# Software implementation

**Super dots analysis tool**
- Java application (available upon request)

**Creating tiles**
- Tiles are 512x512 sized png images (also used by Google Maps, Bing Maps, OSM)
- R package **dotmap**
  - In development: https://github.com/mtennekes/dotmap
  - Both methods (blended colours and super dots) are implemented
  - Working, but no documentation yet

**Visualization**
- R package **tmap** or Javascript library **leaflet**
- Dynamic legend: Javascript

# Acknowledgements



Edwin de Jonge (CBS)



Wouter Meulemans
(TU Eindhoven)



Chantal Melser (CBS)



François Engelen
(Hogeschool Zuyd)