

Stetl: Geo-Conversion-Engine for NLExtract and Smart Emission

www.stetl.org

Just van den Broecke
FOSS4GNL 2018
Almere - The Netherlands
July 11, 2018
www.justobjects.nl



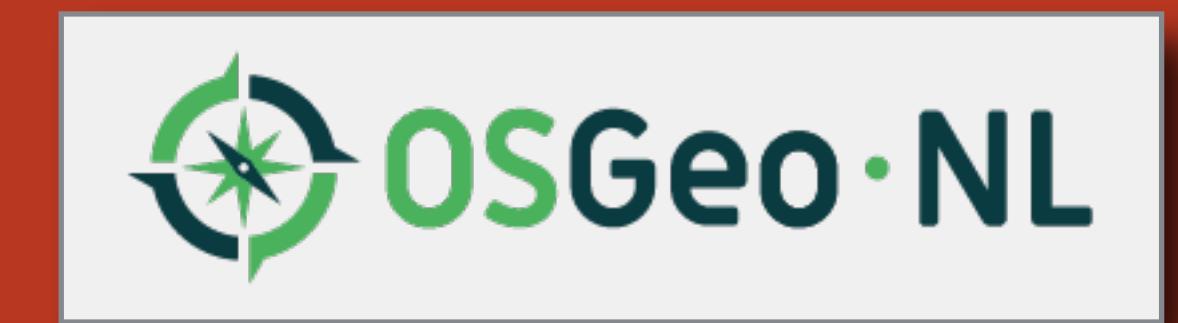
About Me

Independent Open Source Geospatial Professional

- + Chair OSGeo Dutch Local Chapter
- + Member of the Dutch OpenGeoGroep



Just van den Broecke
just@justobjects.nl
www.justobjects.nl

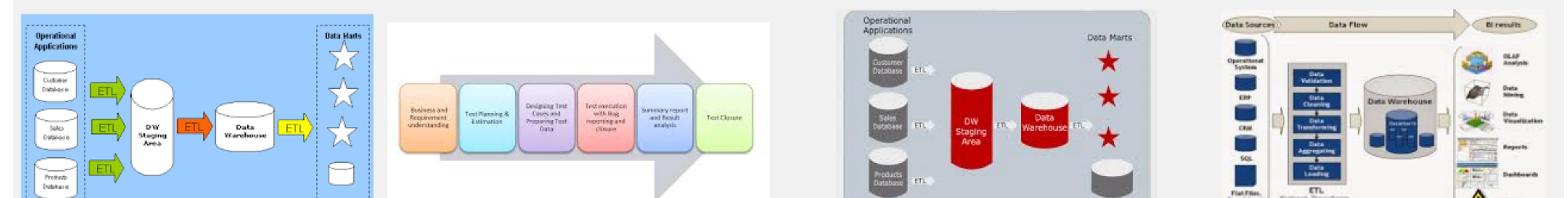
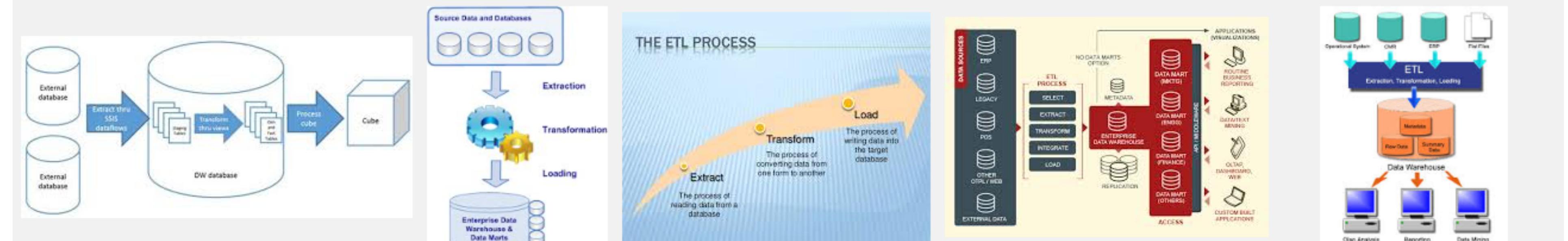
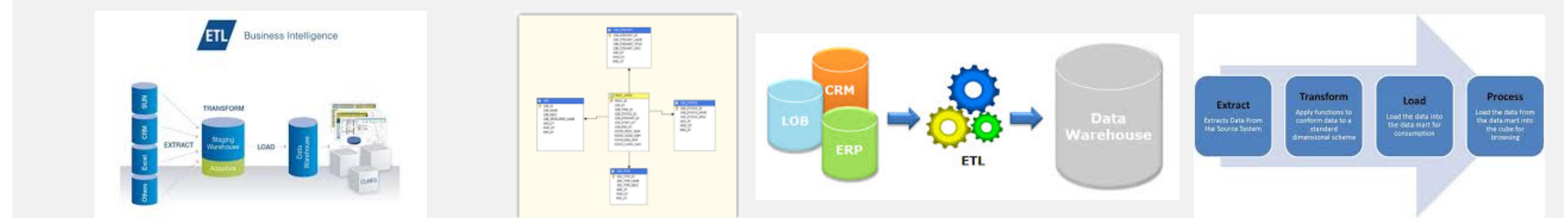
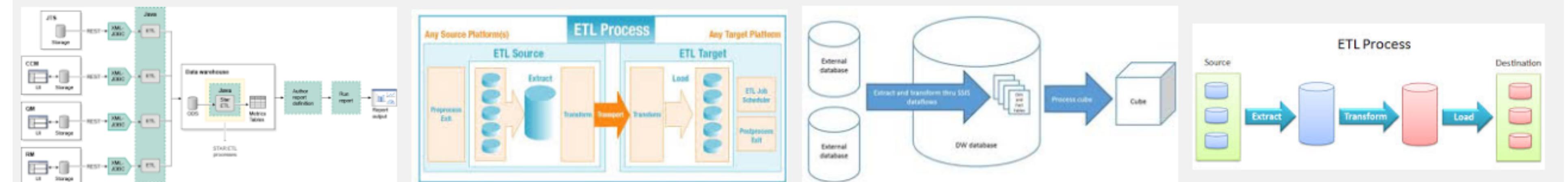
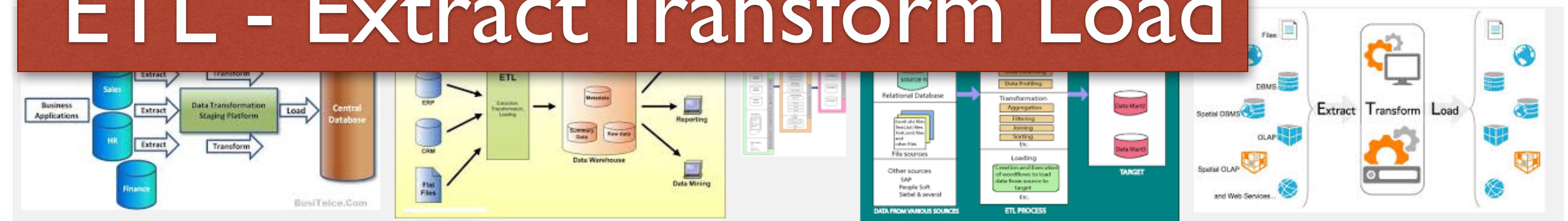


Agenda

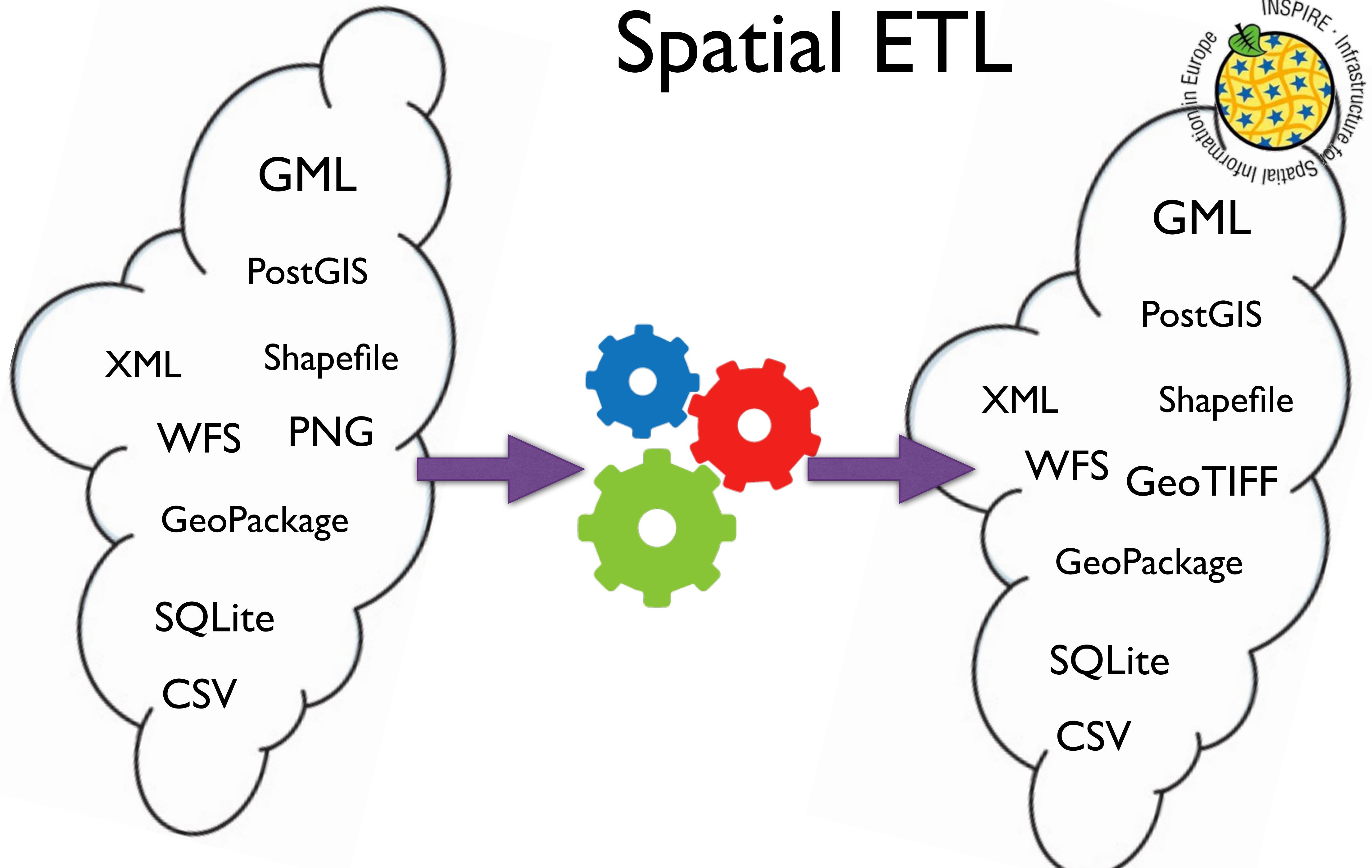
- Spatial ETL
- Stetl Concepts
- Cases
 - NLExtract
 - Smart Emission
- Q & A

Spatial ETL

ETL - Extract Transform Load



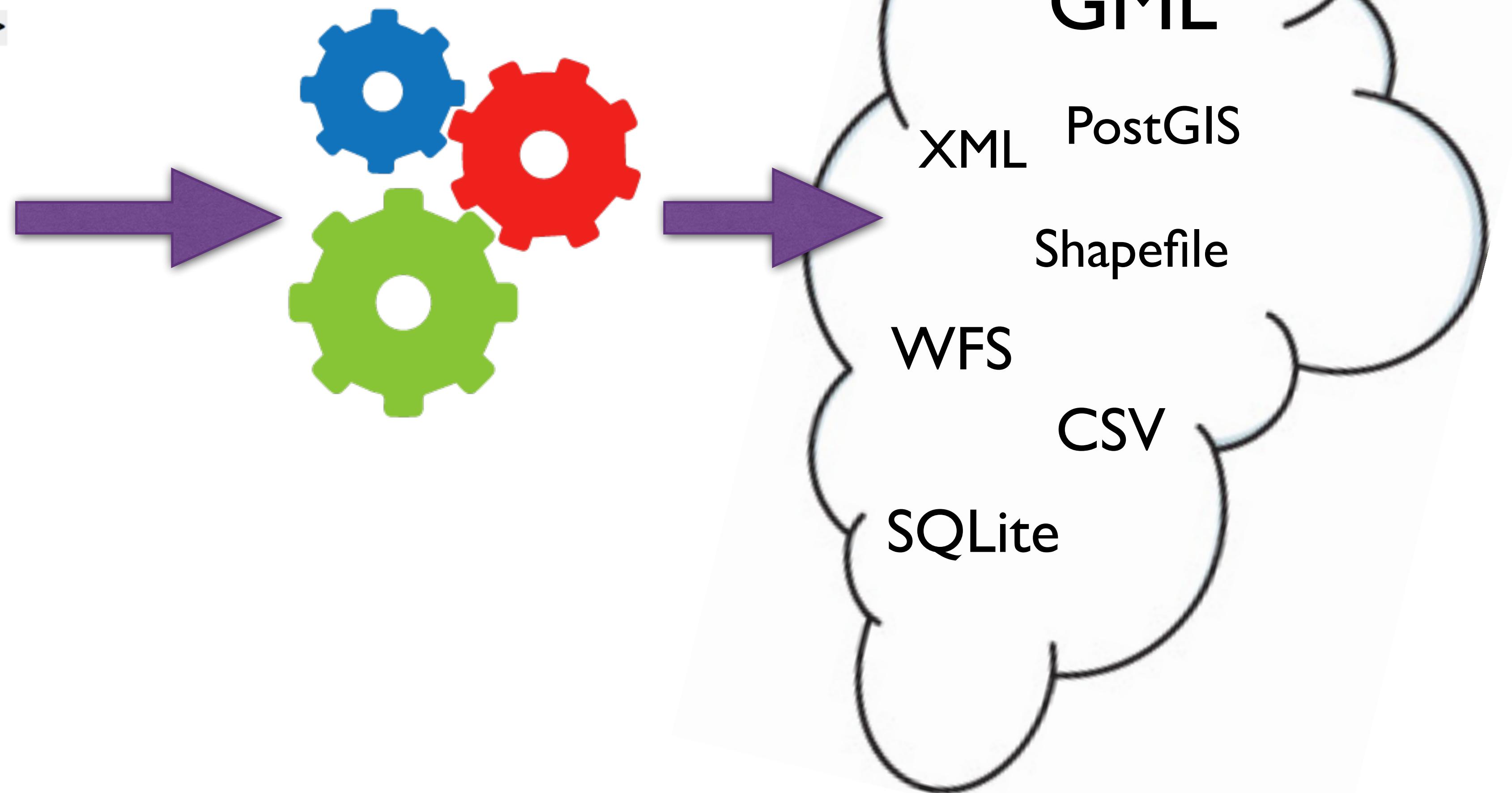
Spatial ETL



Spatial ETL

Example non-standard source

```
<?xml version='1.0' encoding='utf-8'?>
<cities>
  <city>
    <name>Amsterdam</name>
    <lat>52.4</lat>
    <lon>4.9</lon>
  </city>
  <city>
    <name>Bonn</name>
    <lat>50.7</lat>
    <lon>7.1</lon>
  </city>
  <city>
    <name>Rome</name>
    <lat>41.9</lat>
    <lon>12.5</lon>
  </city>
</cities>
```



GDAL/OGR and PROJ.4

<http://www.gdal.org>

<http://proj.osgeo.org>



GDAL for Raster Data - The Geospatial Data Abstraction Library

Arc/Info ASCII Grid, **Arc/Info** Binary Grid (.adf), AIRSAR Polarimetric, Microsoft Windows Device Independent Bitmap (.bmp), BSB Nautical Chart Format (.kap), VTP Binary Terrain Format (.bt), CEOS (Spot for instance), First Generation USGS DOQ (.doq), DODS / OPeNDAP, New Labelled USGS DOQ (.doq), Military Elevation Data (.dt0, .dt1), **ERMapper** Compressed Wavelets (.ecw), **ESRI** .hdr Labelled, ENVI .hdr Labelled Raster, Envisat Image Product (.n1), EOSAT FAST Format, FITS (.fits), Graphics Interchange Format (.gif), GMT Compatible netCDF, **GRASS Rasters**, Golden Software ASCII Grid, Golden Software Binary Grid, Golden Software **Surfer** 7 Binary Grid, TIFF / **GeoTIFF** (.tif), GXF - Grid eXchange File, Hierarchical Data Format Release 4 (**HDF4**), Hierarchical Data Format Release 5 (**HDF5**), **Erdas Imagine** (.img), Vexcel MFF2, **Idrisi** Raster, Image Display and Analysis (WinDisp), ILWIS Raster Map (.mpr,.mpl), Japanese DEM (.mem), **JPEG** JFIF (.jpg), **JPEG2000** (JPEG2000, JP2KAK, JP2ECW, JP2MrSID), NOAA Polar Orbiter Level 1b Data Set (**AVHRR**), Erdas 7.x .LAN and .GIS, Dayton Leveller Heightfield, In Memory Raster, Vexcel MFF, Multi-resolution Seamless Image Database, Meteosat Second Generation, NDF, NITF, **NetCDF**, OGDI Bridge, PCI .aux Labelled, PCI Geomatics Database File, Portable Network Graphics (.png), PCRaster (.map), Netpbm (.ppm,.pgm), Swedish Grid RIK (.rik), RadarSat2 XML (product.xml), **ArcSDE** Raster, USGS SDTS DEM (*CATD.DDF), Raster Matrix Format (*.rsw, .mtw), SAR CEOS, **SGI** Image Format, USGS ASCII DEM (.dem), OGC Web Coverage Server, X11 Pixmap (.xpm)

OGR for Vector Data - Simple Feature Library

Arc/Info Binary Coverage, Comma Separated Value (.csv), DODS/OPeNDAP, DWG, **DXF**, ESRI Personal GeoDatabase, ESRI ArcSDE, **ESRI Shapefile**, FMEObjects Gateway, GML, **GMT Mapping**, **GRASS Vectors**, INTERLIS, **Google Earth KML**, Mapinfo File, Microstation DGN, Spatial **MySQL**, OGDI Vectors, ODBC generic database access layer, **Oracle** Spatial, PostgreSQL **PostGIS**, S-57 (ENC), SDTS, SQLite, UK .NTF, U.S. Census TIGER/Line, VRT - Virtual Datasource, Informix DataBlade

From: https://live.osgeo.org/en/overview/gdal_overview.html



Kersten
@Fernerkundung



Follow

When people describe [#GDAL](#) als the "swiss army knife of GIS data" this is what I think of.



RETWEETS

20

LIKES

21



12:49 PM - 25 Jan 2016



20

21

...

Plenty of Tools...



ogr2ogr

*Each tool is powerful by itself but
cannot do the entire ETL*

FOSS ETL - How to Combine Components?



+

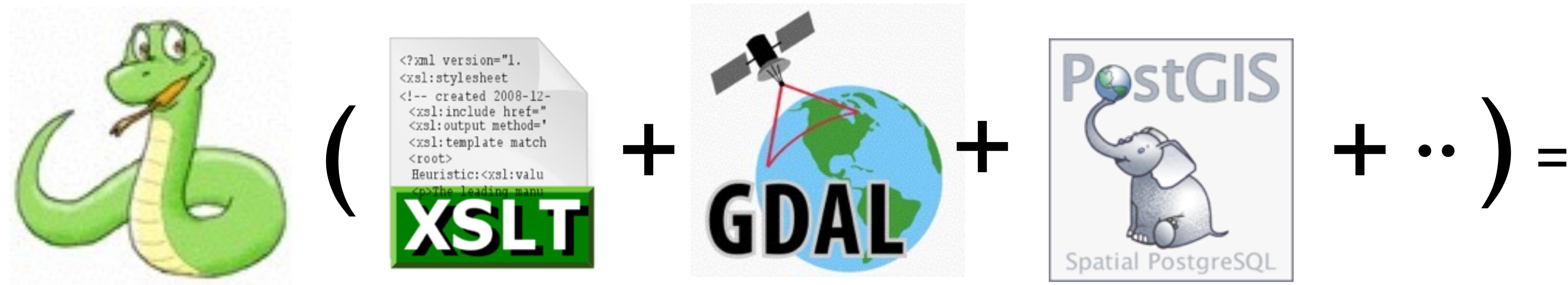


+



+ .. = ?

Solution: Add Python to the Equation



Stetl

Stetl
=

Simple
Streaming
Spatial
Speedy
ETL

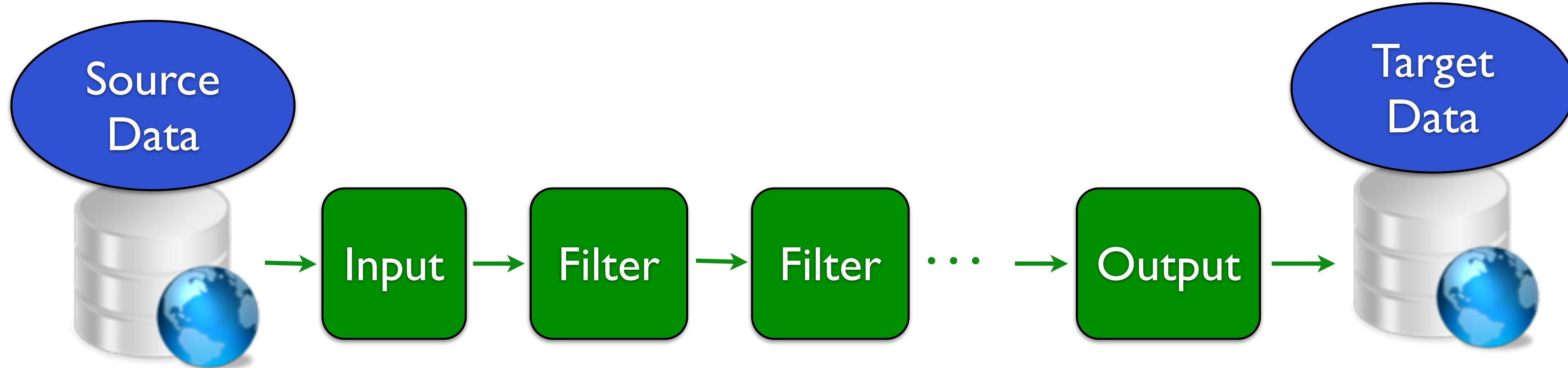


And what about?

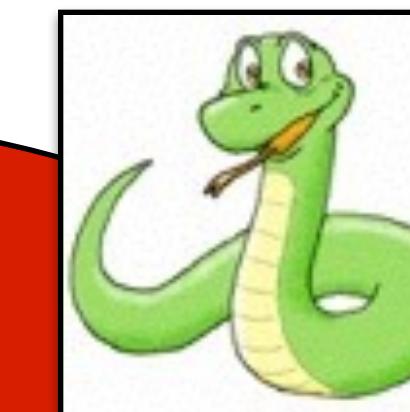


Stetl Concepts

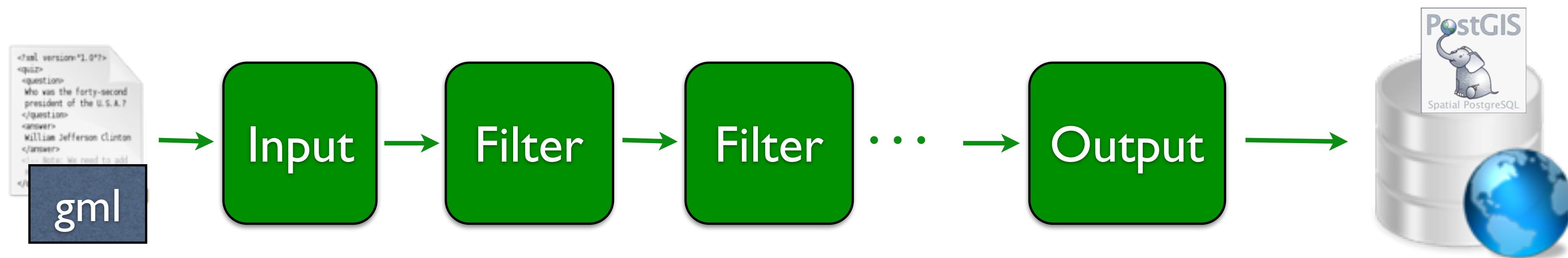
Process Chain



Each Block
is a (reusable) Python Class

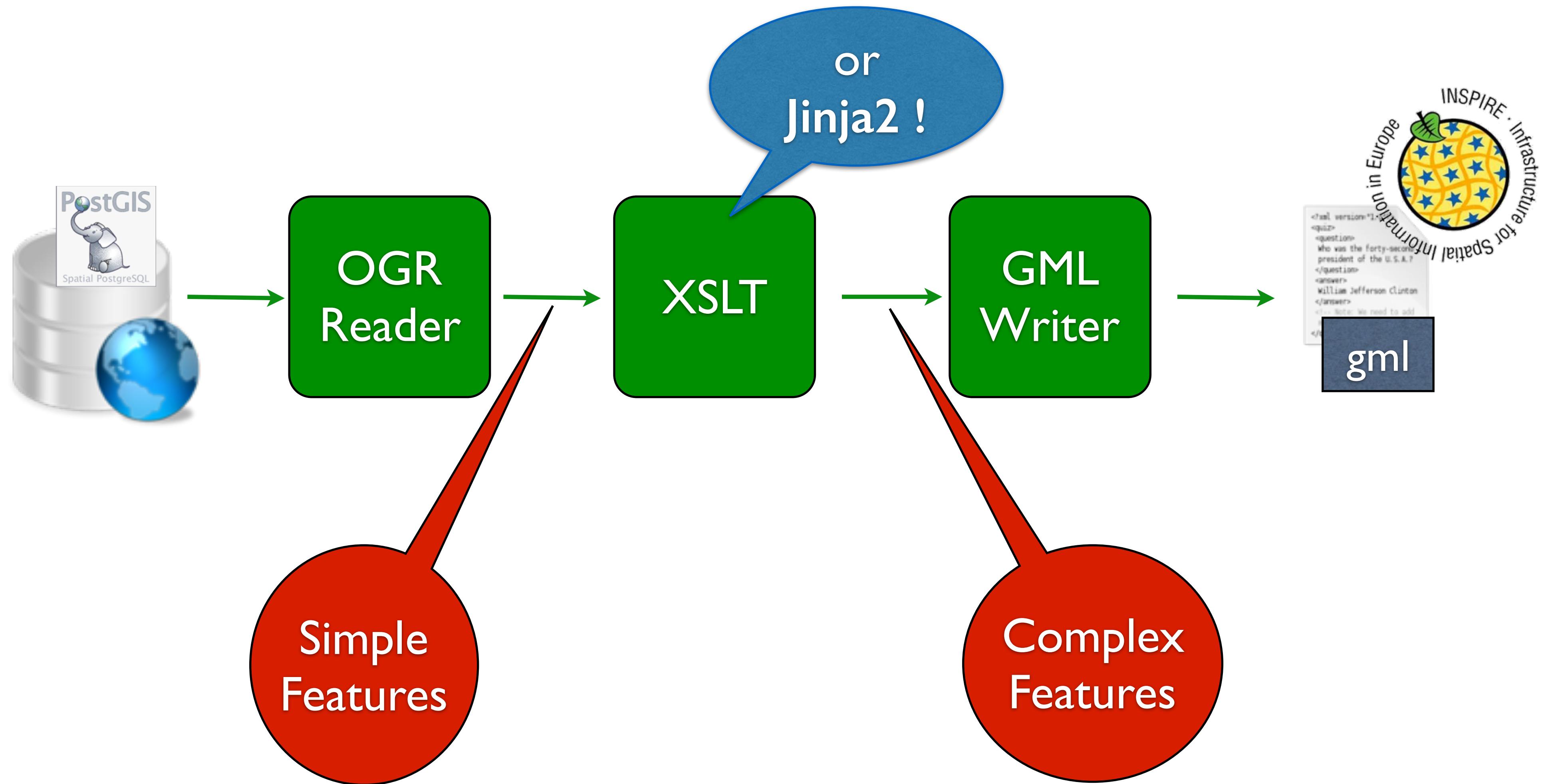


Process Chain - GML to PostGIS



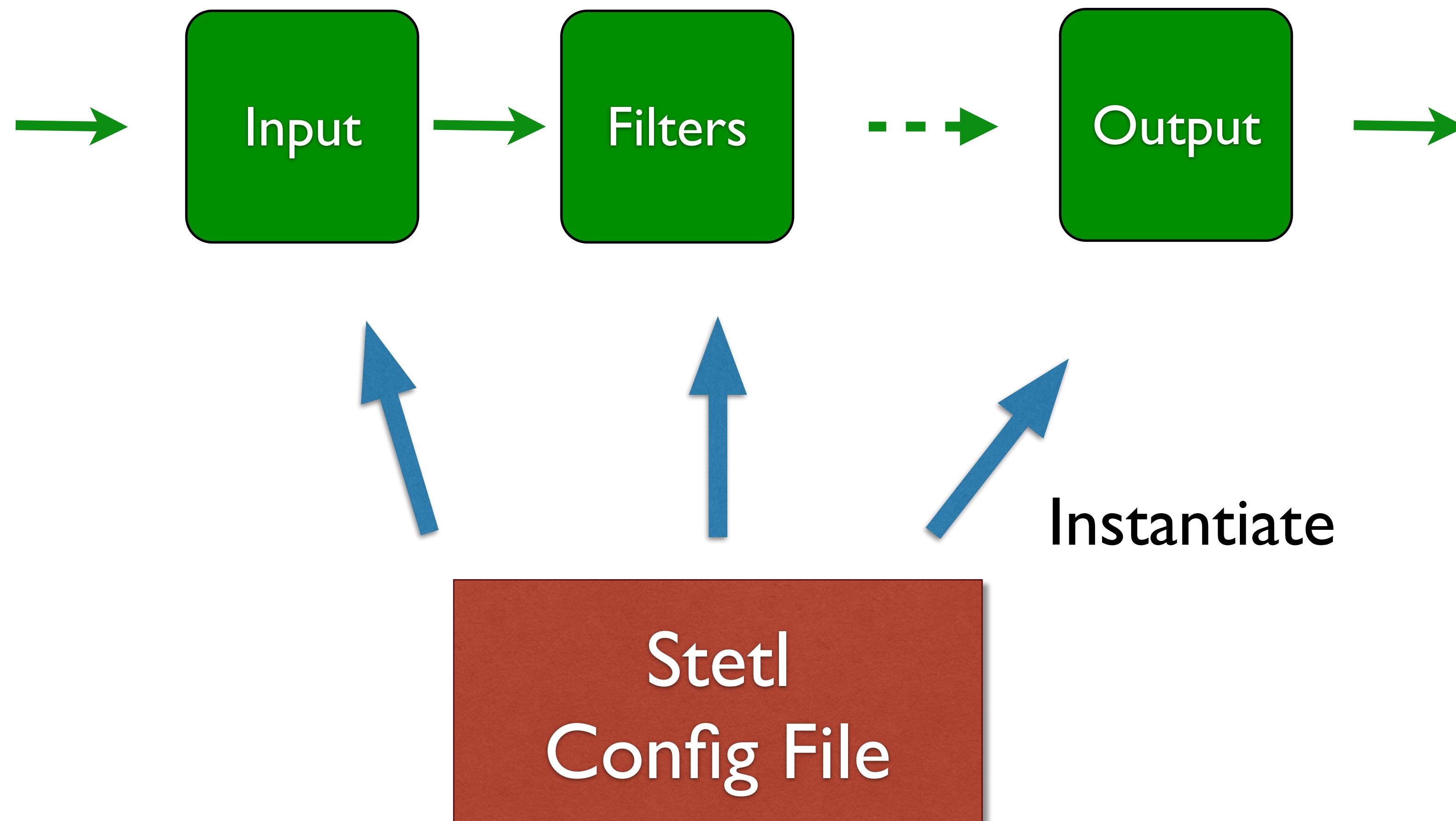
Stetl concepts

Example: Data Model Transform



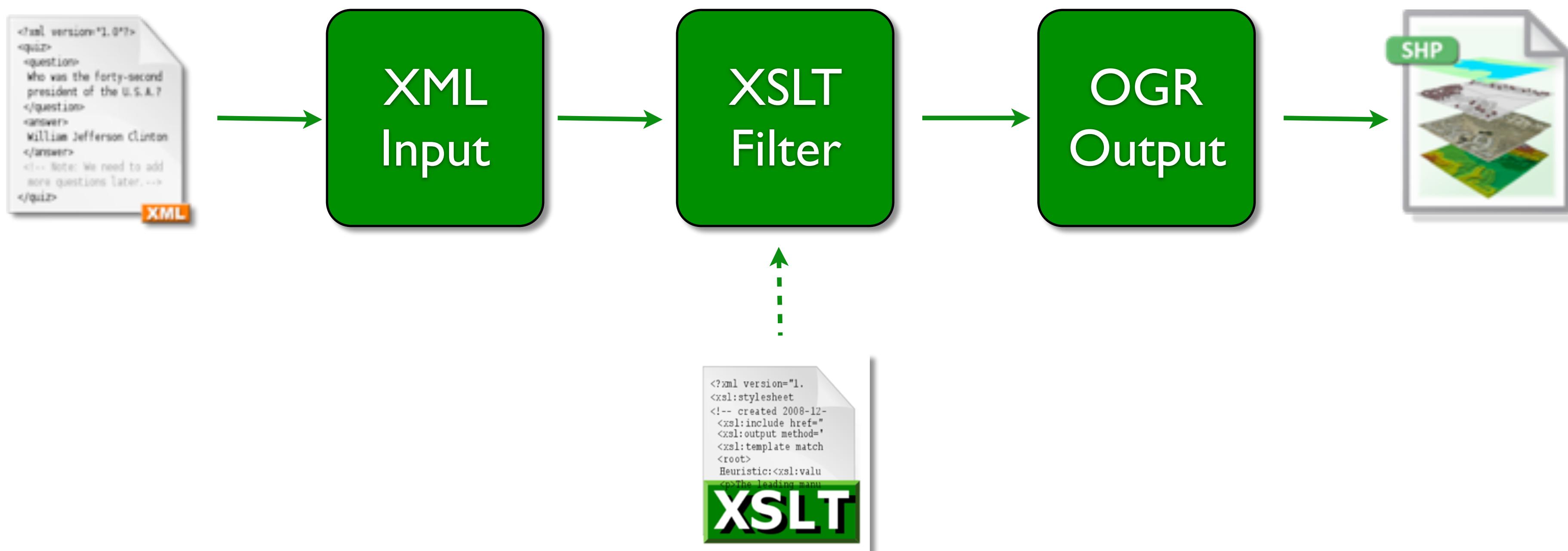
Stetl concepts

Process Chain - How?



Stetl concepts

Example: XML File to Shapefile

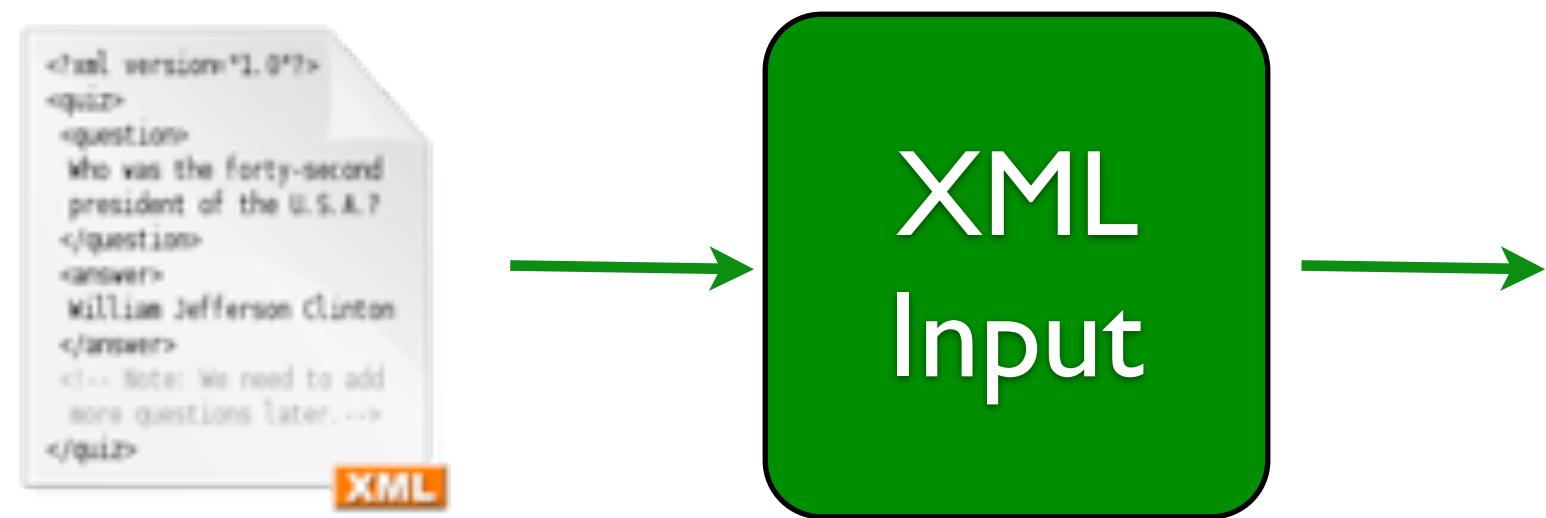


Example: XML to Shape

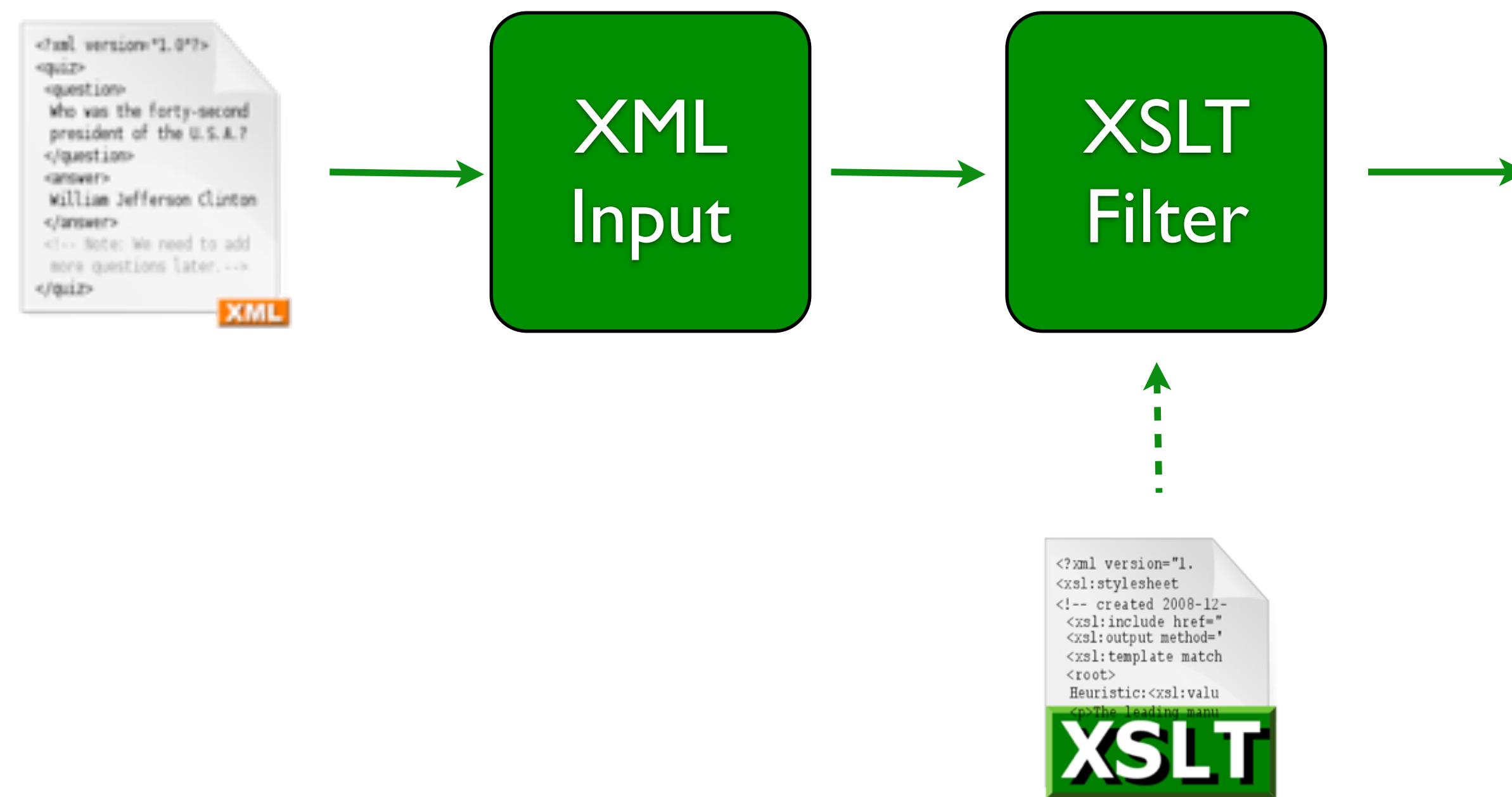
```
1  <?xml version='1.0' encoding='utf-8'?>
2  <cities>
3      <city>
4          <name>Amsterdam</name>
5          <lat>52.4</lat>
6          <lon>4.9</lon>
7      </city>
8      <city>
9          <name>Bonn</name>
10         <lat>50.7</lat>
11         <lon>7.1</lon>
12     </city>
13     <city>
14         <name>Rome</name>
15         <lat>41.9</lat>
16         <lon>12.5</lon>
17     </city>
18 </cities>
```

The Source File

Example: XML to Shape



Example: XML to Shape



Example: XML to Shape

```
<xsl:template match="/">
  <oogr:FeatureCollection
    xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
    xmlns:ogr="http://ogr.maptools.org/"
    xmlns:gml="http://www.opengis.net/gml"
    xsi:schemaLocation="http://ogr.maptools.org/ ../gmlcities.xsd  http://www.opengis.net/gml http://schemas>
  >
  <gml:boundedBy>
    <gml:Box>
      <gml:coord><gml:X>-180.0</gml:X><gml:Y>-90.0</gml:Y></gml:coord>
      <gml:coord><gml:X>180.0</gml:X><gml:Y>90.0</gml:Y></gml:coord>
    </gml:Box>
  </gml:boundedBy>
  <!-- Loop through all cities. -->
  <xsl:apply-templates/>
</ogr:FeatureCollection>
</xsl:template>

<!-- Make each city an ogr:featureMember. -->
<xsl:template match="city">
  <gml:featureMember>
    <ogr:City>
      <ogr:name>
        <xsl:value-of select="name"/>
      </ogr:name>
      <ogr:geometry>
        <gml:Point srsName="urn:ogc:def:crs:EPSG:4326">
          <gml:coordinates><xsl:value-of select="lat"/>,<xsl:value-of select="lon"/></gml:coordinates>
        </gml:Point>
      </ogr:geometry>
    </ogr:City>
  </gml:featureMember>
</xsl:template>
</xsl:stylesheet>
```

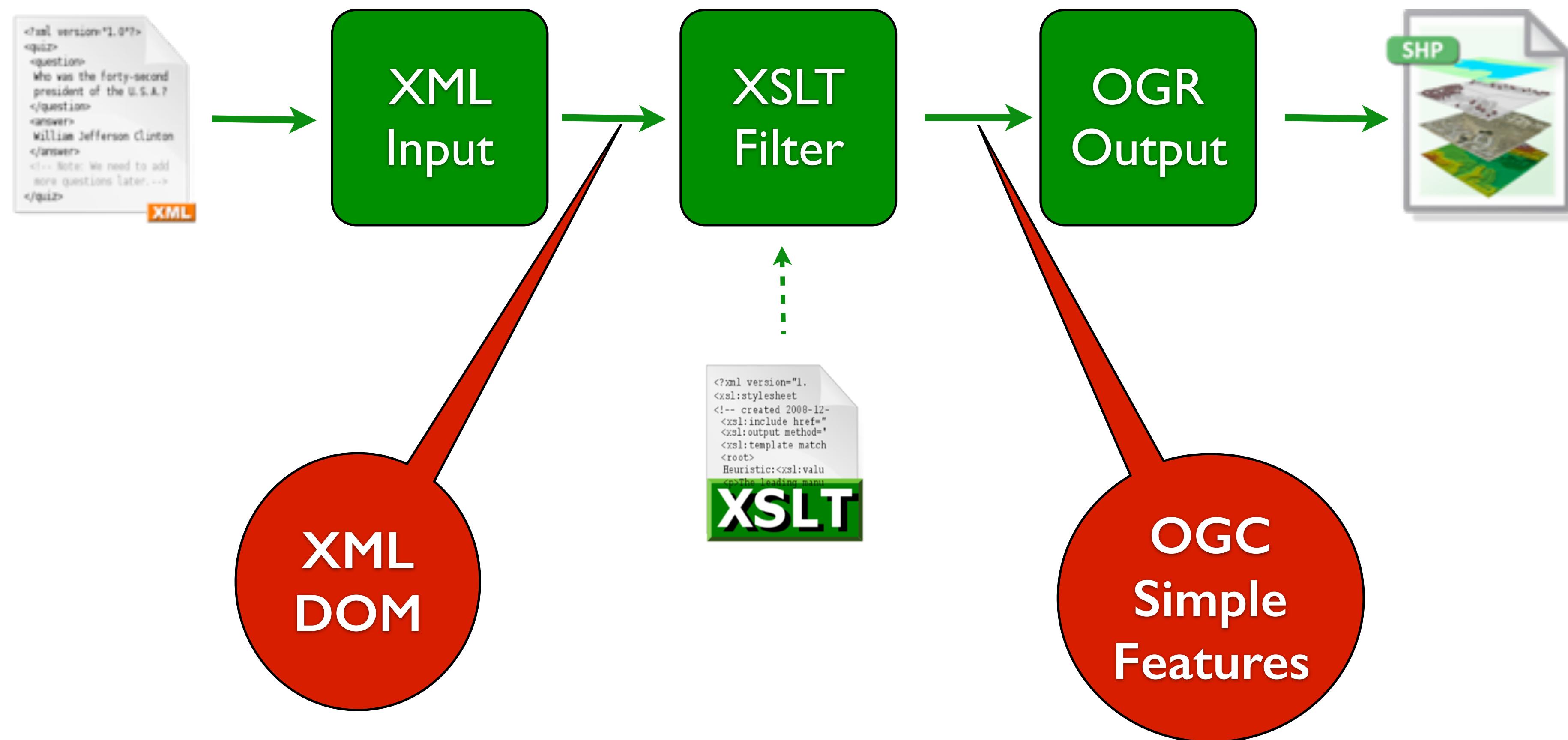
Prepare XSLT Script

Example: XML to Shape

```
1  <?xml version='1.0' encoding='utf-8'?>
2  <oogr:FeatureCollection xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance" xmlns:oogr="http://c
3  <gml:boundedBy>
4  <gml:Box>
5  <gml:coord>
6  <gml:X>-180.0</gml:X>
7  <gml:Y>-90.0</gml:Y>
8  </gml:coord>
9  <gml:coord>
10 <gml:X>180.0</gml:X>
11 <gml:Y>90.0</gml:Y>
12 </gml:coord>
13 </gml:Box>
14 </gml:boundedBy>
15 <gml:featureMember>
16 <oogr:City>
17 <oogr:name>Amsterdam</oogr:name>
18 <oogr:geometry>
19 <gml:Point srsName="urn:ogc:def:crs:EPSG:4326">
20 <gml:coordinates>52.4,4.9</gml:coordinates>
21 </gml:Point>
22 </oogr:geometry>
23 </oogr:City>
24 </gml:featureMember>
25 <gml:featureMember>
26 <oogr:City>
27 <oogr:name>Bonn</oogr:name>
28 <oogr:geometry>
29 <gml:Point srsName="urn:ogc:def:crs:EPSG:4326">
30 <gml:coordinates>50.7,7.1</gml:coordinates>
31 </gml:Point>
32 </oogr:geometry>
33 </oogr:City>
34 </gml:featureMember>
35 <gml:featureMember>
```

XSLT GML Output

Example: XML to Shape



Example: XML to Shape

```
# Transform input xml to valid GML file and then to Shape.

[etl]
chains = input_xml_file | transformer_xslt | output_ogr_shape

[input_xml_file]
class = inputs.fileinput.XmlFileInput
file_path = input/cities.xml

[transformer_xslt]
class = filters.xsltfilter.XsltFilter
script = cities2gml.xsl

# The ogr2ogr command-line.
[output_ogr_shape]
class = outputs.ogroutput.Ogr2OgrOutput
temp_file = temp/gmlcities.gml
ogr2ogr_cmd = ogr2ogr \
    -overwrite \
    -f "ESRI Shapefile" \
    -a_srs epsg:4326 \
    output/gmlcities.shp \
    temp/gmlcities.gml
```

XSLT
Filter

Process
Chain

XML
Input

ogr2ogr
Output

The Stetl Config File

Data Structures

- Components exchange *Packets*
- *Packet* contains data and status
- Data formats, e.g.:

xml_line_stream

etree_doc

etree_element (feature)

etree_element_array

record (dict)

string

any

.

.

Running Stetl

```
stetl -c etl.cfg
```

```
stetl -c etl.cfg [-a <properties>]
```

Example Components



XMLFile

ogr2ogr

LineStream

Postgres/PostGIS

CSV, JSON etc

YourInput

XSLT

XMLAssembler

XMLValidator

Jinja2

FeatureExtractor

YourFilter

GML

GDAL/OGR

WFS-T

Postgres/PostGIS

SOS, SensorThings API

YourOutput

Example: XsltFilter Python Code

```
class XsltFilter(Filter):
    """
    Invokes XSLT processor (via lxml) for given XSLT script on an etree doc.
    """

    @Config(ptype=str, required=True)
    def script(self):
        """
        Path to XSLT script file.
        """
        pass

    def __init__(self, configdict, section):
        Filter.__init__(self, configdict, section, consumes=FORMAT.etree_doc, produces=FORMAT.etree_doc)

        self.xslt_file = open(self.script, 'r')

        # Parse XSLT file only once
        self.xslt_doc = etree.parse(self.xslt_file)
        self.xslt_obj = etree.XSLT(self.xslt_doc)
        self.xslt_file.close()

    def invoke(self, packet):
        packet.data = self.xslt_obj(packet.data)
        return packet
```

Your Own Components

```
class MyFilter(Filter):
    # Constructor
    def __init__(self, configdict, section):
        Filter.__init__(self, configdict, section, consumes=FORMAT.etree_doc,
                       produces=FORMAT.etree_doc)

    def invoke(self, packet):
        log.info("CALLING MyFilter OK!!!!")
        return packet
```

Step 1 - Define Class

```
[etl]
chains = input_xml_file|my_filter|output_std

[input_xml_file]
class = inputs.fileinput.XmlFileInput
file_path = input/cities.xml

# My custom component
[my_filter]
class = my.myfilter.MyFilter

[output_std]
class = outputs.standardoutput.StandardXmlOutput
```

Step 2- Configure Class

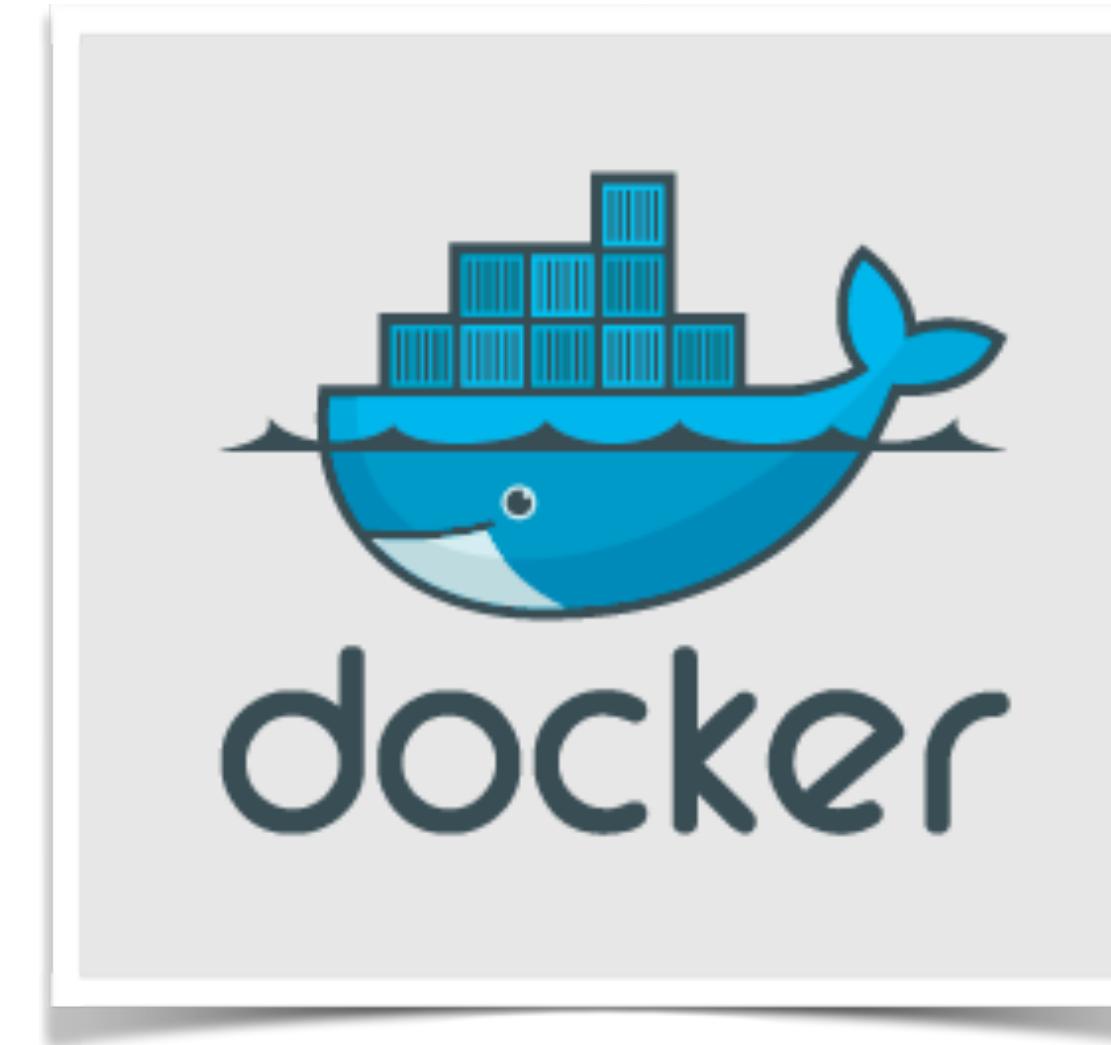
Installing Stetl - PyPi

```
sudo pip install stetl
```

Deps

- GDAL+Python bindings
- lxml (xml proc)
- psycopg2 (Postgres)

Installing Stetl - Docker



[https://hub.docker.com/r/
geopython/stetl](https://hub.docker.com/r/geopython/stetl)

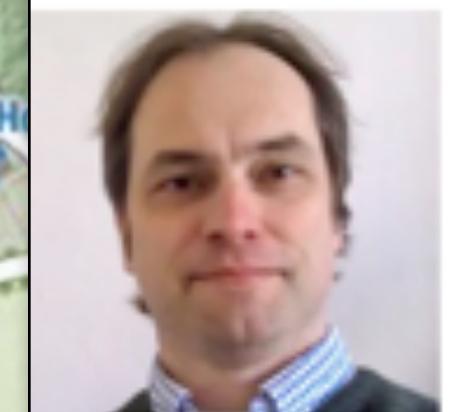
Case I - NLExtract

What is NLExtract?

ETL for Dutch National Open Datasets

<http://nlextract.nl>

- GML to PostGIS:
BAG, BRT (Top10NL), BGT, BRK (Cadastral Parcels)
- Downloads data.nlextract.nl

**Jan-Willem van Aalst**Imergis organisatiebloei
Utrecht Area, Netherlands | Professional TraiCurrent [www.imergis.nl](#)
Previous [GHD Nederland, Ministerie VenJ,](#)
Education [Technische Universiteit Delft](#)[Send a message](#)

Otterlo

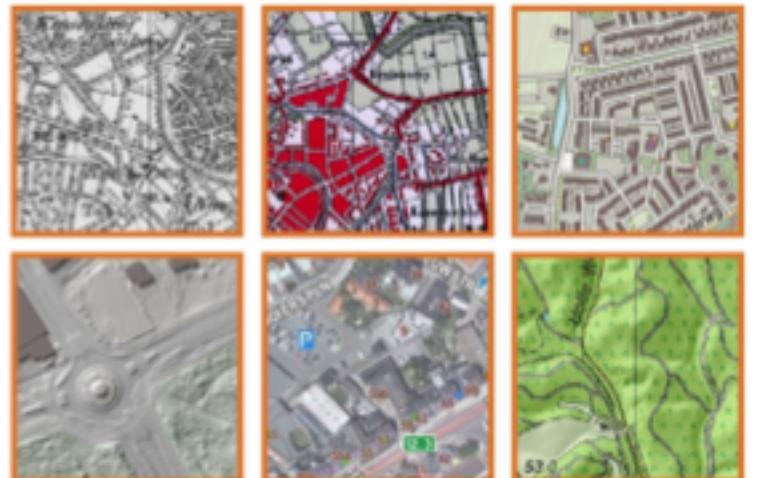


Map5.nl

Map5.nl: de fijnste kaarten voor al je toepassingen

Topografische kaarten van Nederland via open geo webstandaarden

- Topografische kaarten: OpenTopo, OpenSimpleTopo, Kadaster Top250/50/25
- Luchtfoto's met OpenTopo labels
- Reliëfkaarten op basis AHN2 (0.5m resolutie)
- Historische kaarten: Bonnebladen, TMK 1850
- Hoge schalen (tot 1:150)
- Standaarden: TMS, WMTS, WMS en Google/OSM x,y,z
- Tiling schema's: Nederlands (RD/PDOK) en Web Mercator (Google/OSM)
- Gemakkelijk opnemen in toepassingen: web, desktop, mobiel
- Mobiele apps
- Gebouwd met Open Source geo-componenten
- CORS headers voor blyv 3D rendering
- Gratis of betaalde dienst
- Eigen unieke URL: geen whitelisting, certificaten of API-keys



Bekijk in de NLTopo App, ook op je mobiel!

Wat kost het?

BRK - Input GML from PDOK

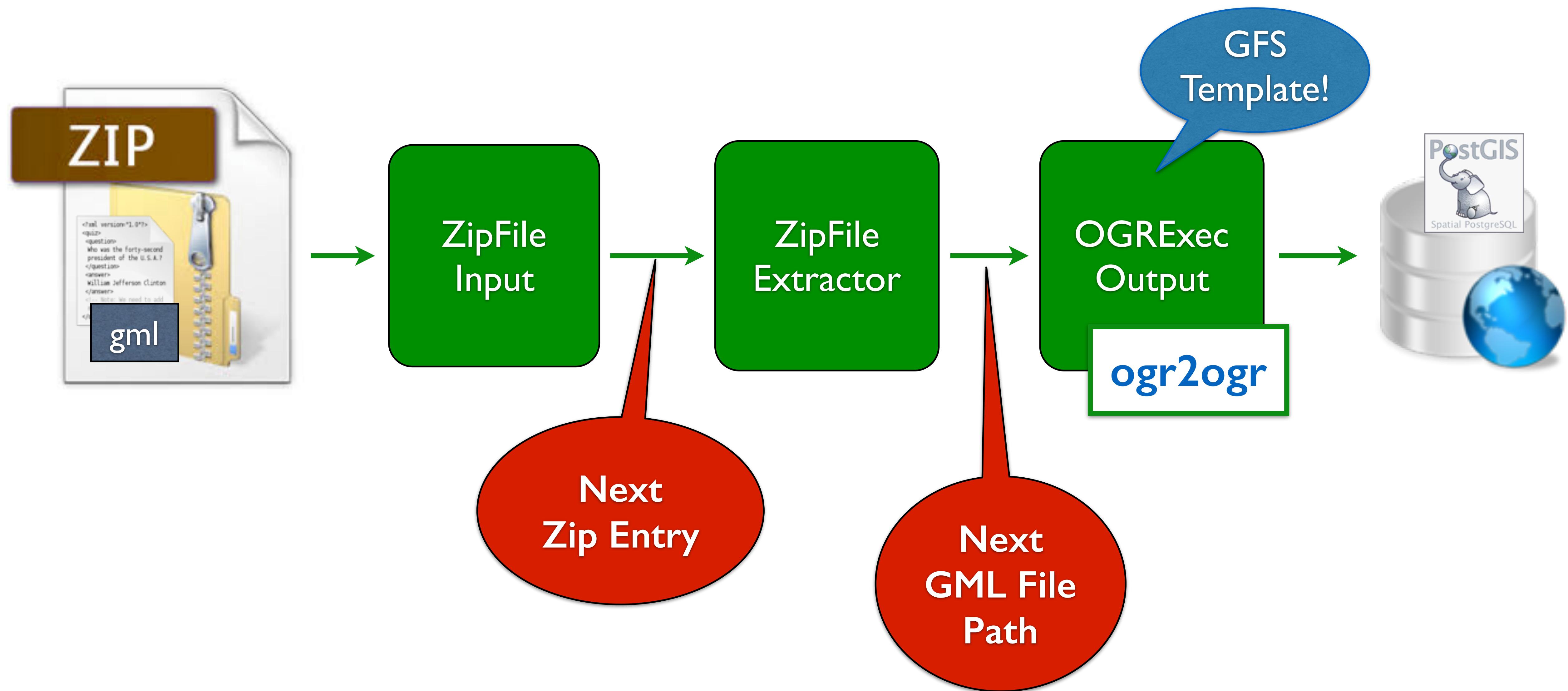
```
$ ls -lh  
total 3.6G  
-rw-rw-r-- 1 mlx mlx 3.6G May 7 10:32 kadastralekaartv3-gml-nl-nohist.zip
```

```
$ unzip -l kadastralekaartv3-gml-nl-nohist.zip
```

Length	Date	Time	Name
-----	-----	-----	-----
16562827212	2018-01-01	13:42	Bebouwing.gml
4737201554	2018-01-02	09:00	Annotatie.gml
20883080403	2018-05-01	15:15	Perceel.gml
27532931147	2018-05-01	20:19	Kadastralegrens.gml
-----			-----
69716040316			4 files

About 70GB in 4 GML Files!

Stetl Config BRK to PostGIS



Stetl Config - BRK .zip to PostGIS

```
[etl]
chains = input_zip_file|extract_zip_file|output_ogr2ogr

[input_zip_file]
class=inputs.fileinput.ZipFileInput
file_path = {input_dir}
filename_pattern = {zip_files_pattern}
name_filter = {filename_match}

# Filter to extract a ZIP file one by one to a temporary location
[extract_zip_file]
class=filters.zipfileextractor.ZipFileExtractor
file_path = {temp_dir}/fromzip-tmp.gml

[output_ogr2ogr]
class = outputs.execoutput.Ogr2OgrExecOutput
# destination format: OGR vector format name
dest_format = PostgreSQL
# destination datasource: name of datasource
dest_data_source = "PG:dbname={database} host={host} port={port} user={user} password={password} active_schema={schema}"
# layer creation options will only be added to ogr2ogr on first run
lco = -lco LAUNDER=YES -lco PRECISION=NO
# spatial_extent, translates to -spat xmin ymin xmax ymax
spatial_extent = {spatial_extent}
# gfs template
gfs_template = {gfs_template}
# miscellaneous ogr2ogr options
options = -append -gt 65536 {multi_opts} --config PG_USE_COPY NO --config CPL_ZIP_ENCODING CP437
```

GFS for Table+Column Mapping

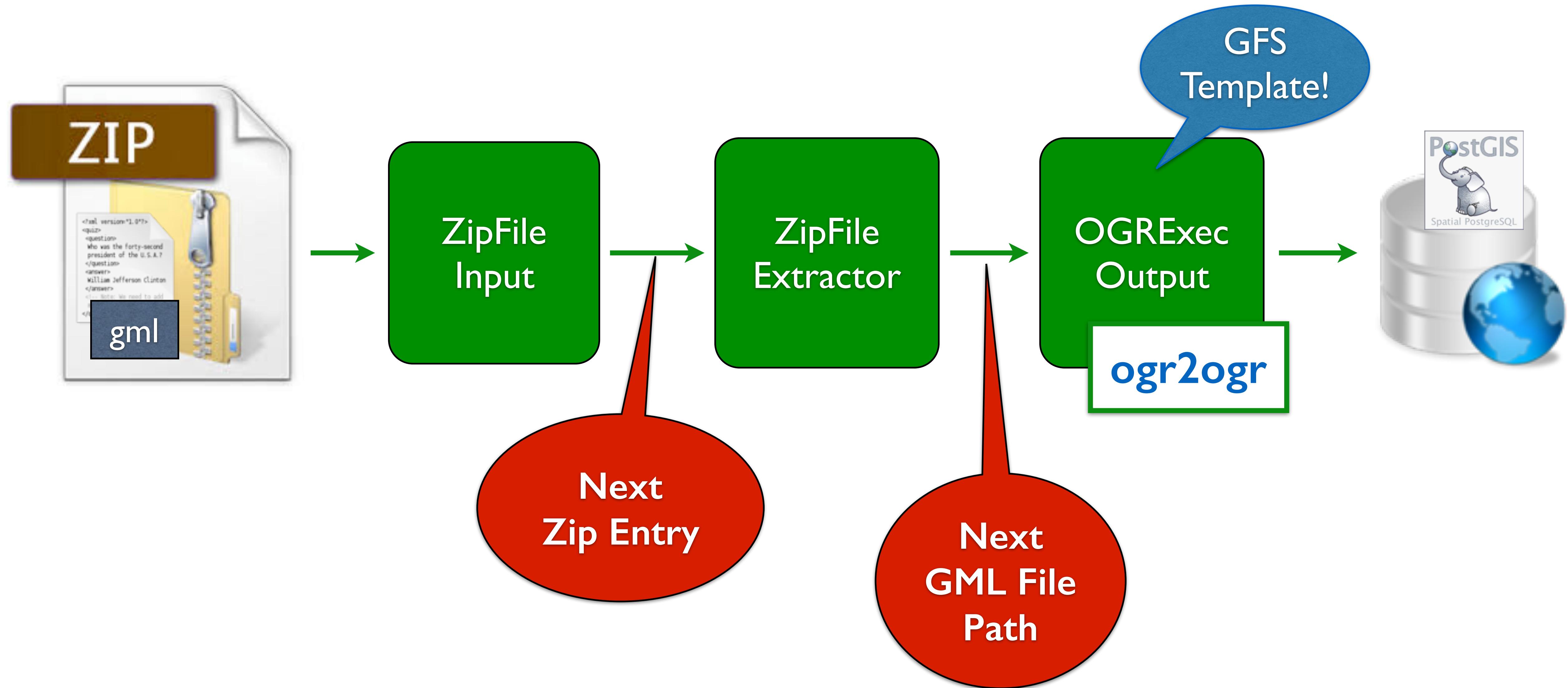
```
<GMLFeatureclassList>
  .
  <GMLFeatureClass>
    <Name>Perceel</Name>
    <ElementPath>Perceel</ElementPath>
    <!--<GeometryType>1</GeometryType>-->
    <SRSName>urn:ogc:def:crs:EPSG::28992</SRSName>
    .
    <PropertyDefn>
      <Name>gemeente</Name>
      <ElementPath>kadastraleAanduiding|TypeKadastraleAanduiding|kadastraleGemeente|KadastraleGemeenteKeuze|
AKRKadastraleGemeenteCode</ElementPath>
      <Type>String</Type>
      <Width>5</Width>
    </PropertyDefn>
    <PropertyDefn>
      <Name>sectie</Name>
      <ElementPath>kadastraleAanduiding|TypeKadastraleAanduiding|sectie</ElementPath>
      <Type>String</Type>
      <Width>2</Width>
    </PropertyDefn>
    <PropertyDefn>
      <Name>perceelnummer</Name>
      <ElementPath>kadastraleAanduiding|TypeKadastraleAanduiding|perceelnummer</ElementPath>
      <Type>Integer</Type>
    </PropertyDefn>
    .
    <GeomPropertyDefn>
      <Name>begrenzing</Name>
      <!-- OGR geometry name -->
      <ElementPath>begrenzingPerceel</ElementPath>
      <!-- XML element name possibly with '/' to specify the path -->
      <Type>Polygon</Type>
    </GeomPropertyDefn>
  </GMLFeatureClass>
</GMLFeatureclassList>
```

Resultaat in PostGIS

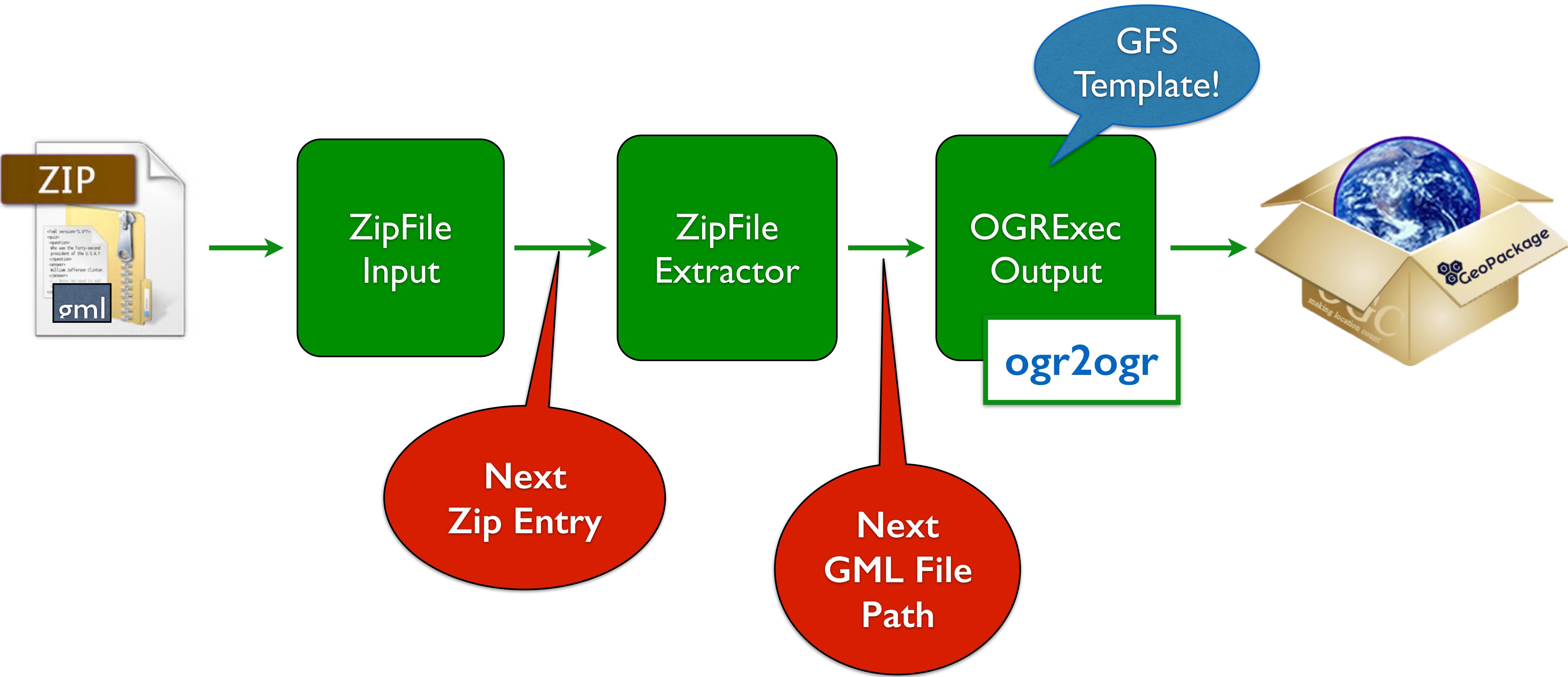
ogc_fid	gml_id	namespace	lokaalid	logischtijdstipontstaan	gemeente	sectie	perceelnummer	waarde
1	NL.IMKAD.KadastraalObject.120031707	NL.IMKAD.KadastraalObject	120031707	2009-03-27T23:59:59.000	AMR04	O	3322	178
2	NL.IMKAD.KadastraalObject.120118807	NL.IMKAD.KadastraalObject	120118807	2009-03-27T23:59:59.000	AMR04	P	3907	446
3	NL.IMKAD.KadastraalObject.120077624	NL.IMKAD.KadastraalObject	120077624	2009-03-27T23:59:59.000	AMR04	O	4482	298
4	NL.IMKAD.KadastraalObject.120119310	NL.IMKAD.KadastraalObject	120119310	2009-03-27T23:59:59.000	AMR04	Q	5520	155
5	NL.IMKAD.KadastraalObject.120095784	NL.IMKAD.KadastraalObject	120095784	2009-03-27T23:59:59.000	AMR04	Q	4128	250

perceelnummer	rotatie	deltax	deltay	begrenzing
	0	0	0	010300002040710000010000000C0000003D0AD7A3CF5E014...
	0	0	0	01030000204071000001000000060000006F1283C072D0014...
	0	0	0	0103000020407100000100000006000000D34D6210DB6F014...
-45.6	0	0	0	010300002040710000010000000900000046B6F3FD0C2D014...
	0	0	0	010300002040710000010000000A000000C74B3789BE42014...

Stetl Config BRK to PostGIS



Stetl Config BRK to GeoPackage



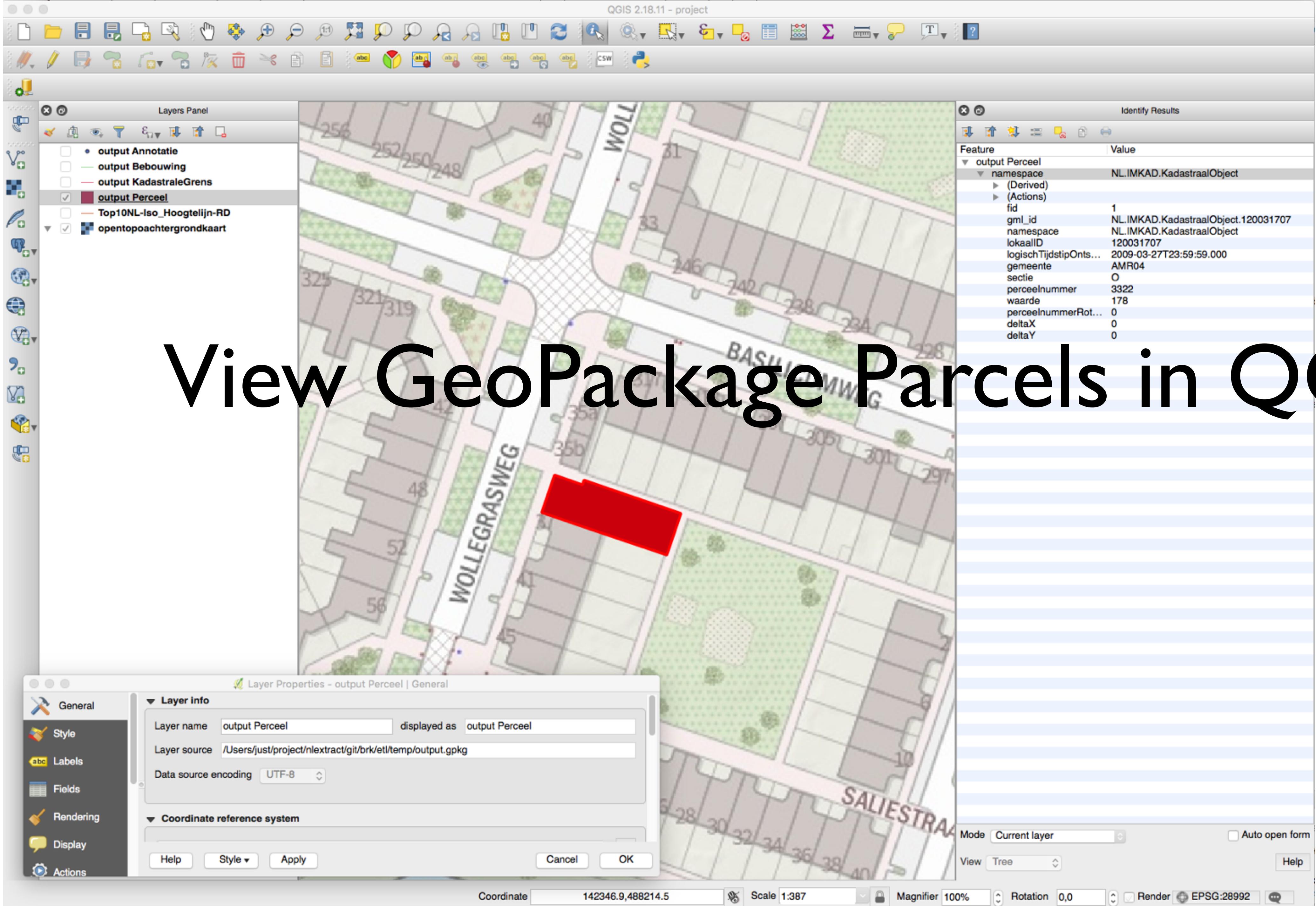
Stetl Config - BRK .zip to GeoPackage

```
[etl]
chains = input_zip_file|extract_zip_file|output_ogr2ogr_gpkg

# The source input ZIP-file(s) from dir, producing 'records' with ZIP file name and inner file names
[input_zip_file]
class=inputs.fileinput.ZipFileInput
file_path = {input_dir}
filename_pattern = {zip_files_pattern}
name_filter = {filename_match}

# Filter to extract a ZIP file one by one to a temporary location
[extract_zip_file]
class=filters.zipfileextractor.ZipFileExtractor
file_path = {temp_dir}/fromzip-tmp.gml

[output_ogr2ogr_gpkg]
class = outputs.execoutput.Ogr2OgrExecOutput
# destination format: OGR vector format name
dest_format = GPKG
# destination datasource: path to .gpkg output file
dest_data_source = "{temp_dir}/output.gpkg"
# layer creation options will only be added to ogr2ogr on first run
lco = -lco SPATIAL_INDEX=YES -lco PRECISION=NO
# spatial_extent, translates to -spat xmin ymin xmax ymax
spatial_extent = {spatial_extent}
# gfs template
gfs_template = {gfs_template}
# miscellaneous ogr2ogr options
options = -append -gt 65536 {multi_opts} --config CPL_ZIP_ENCODING CP437
```



View GeoPackage Parcels in QGIS

BRK .zip to both PostGIS + GeoPackage

```
[etl]
chains = input_zip_file|extract_zip_file|(output_ogr2ogr_postgis) (output_ogr2ogr_gpkg)

# The source input ZIP-file(s) from dir, producing 'records' with ZIP file name and inner file names
[input_zip_file]
class=inputs.fileinput.ZipFileInput
file_path = {input_dir}
filename_pattern = {zip_files_pattern}
name_filter = {filename_match}

# Filter to extract a ZIP file one by one to a temporary location
[extract_zip_file]
class=filters.zipfileextractor.ZipFileExtractor
file_path = {temp_dir}/fromzip-tmp.gml

[output_ogr2ogr_postgis]
class = outputs.execoutput.Ogr2OgrExecOutput
# destination format: OGR vector format name
dest_format = PostgreSQL
.

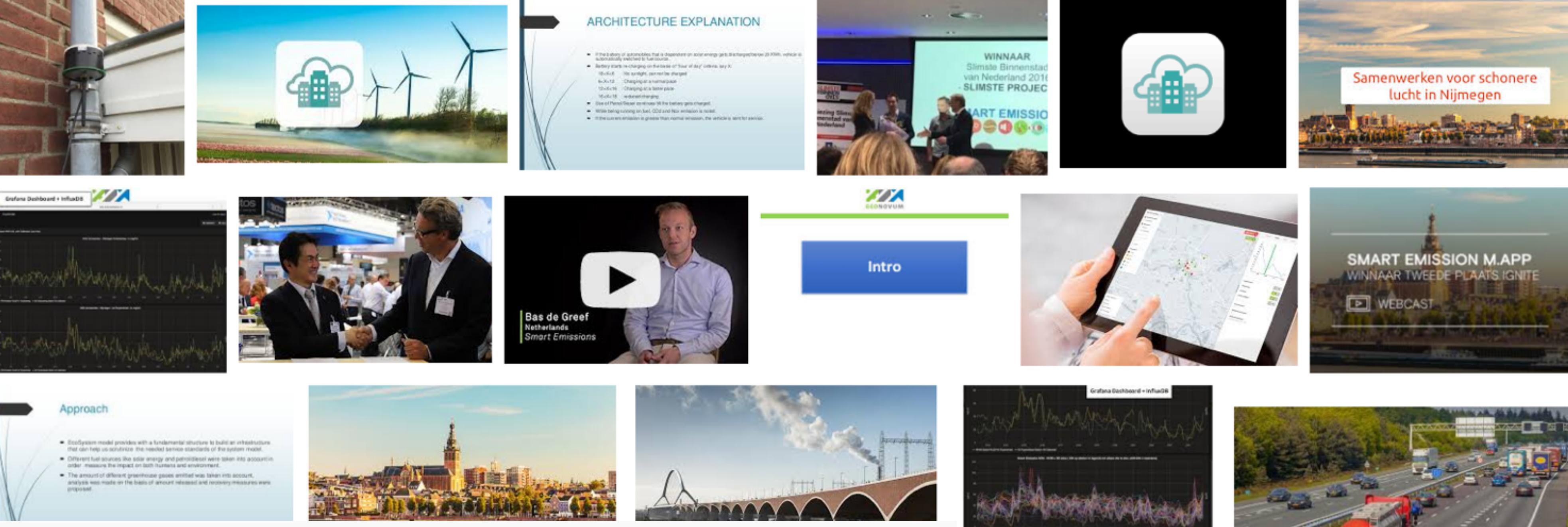
.

[output_ogr2ogr_gpkg]
class = outputs.execoutput.Ogr2OgrExecOutput
# destination format: OGR vector format name
dest_format = GPKG
# destination datasource: path to .gpkg output file
dest_data_source = "{temp_dir}/output.gpkg"
.
```

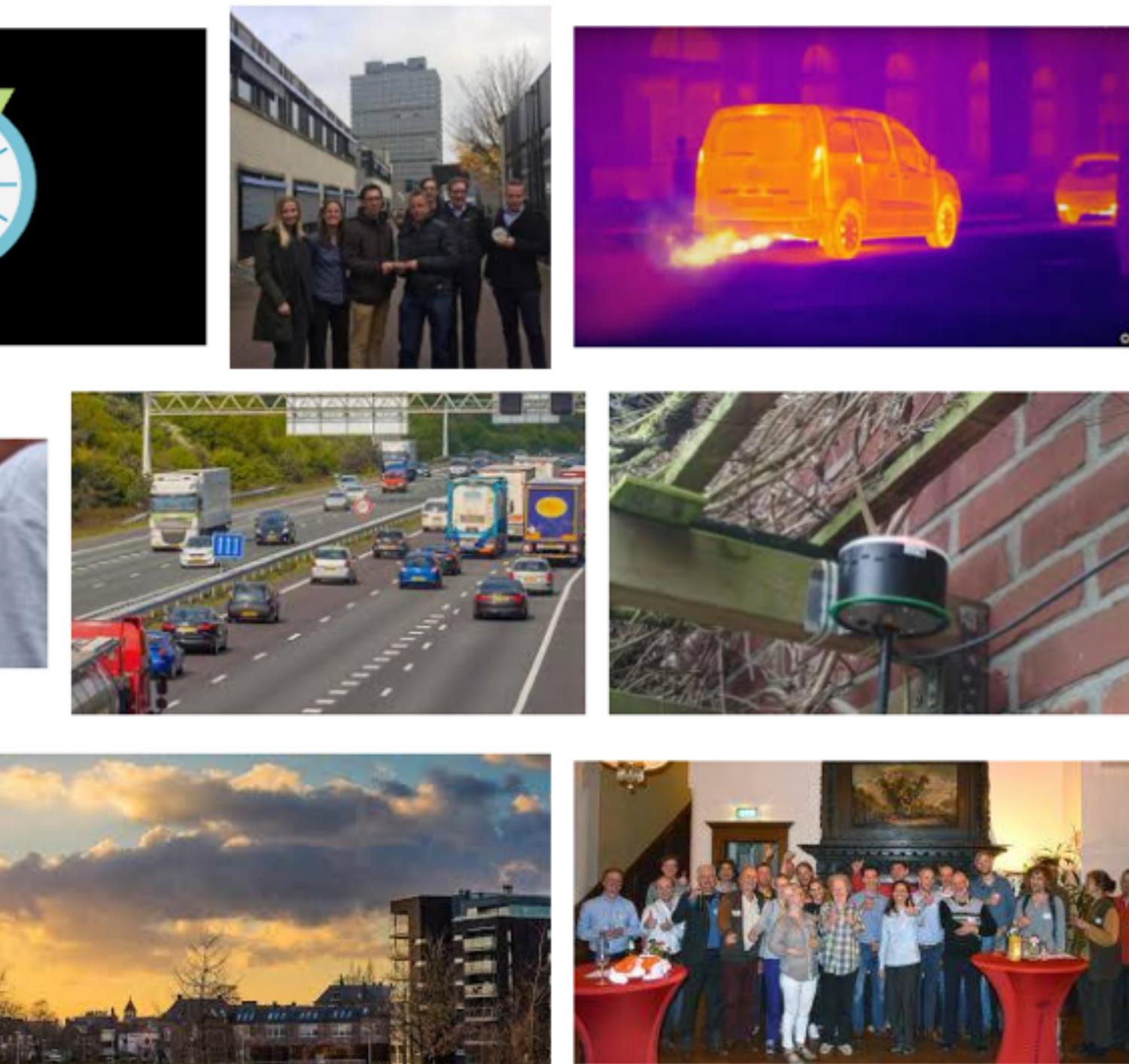


Embeds
Stel
“Splitter”

Case 2 - Smart Emission



Het 'Smart Emission' project draait om het in kaart brengen van luchtkwaliteit, geluid, trillingen en meteorologische indicatoren in de stad op een fijnmazig schaalniveau, door inwoners met zogenoemde burger-sensor-netwerken. De pilotstudie richt zich op de stad Nijmegen (binnen de gemeentegrenzen van de stad). Hier zal een testbed worden ingericht waar het projectteam, in samenwerking met inwoners de lokale variatie in luchtkwaliteit (NO₂, CO, CO₂, O₃), geluid en trillingen in kaart brengen met behulp van betaalbare sensoren.

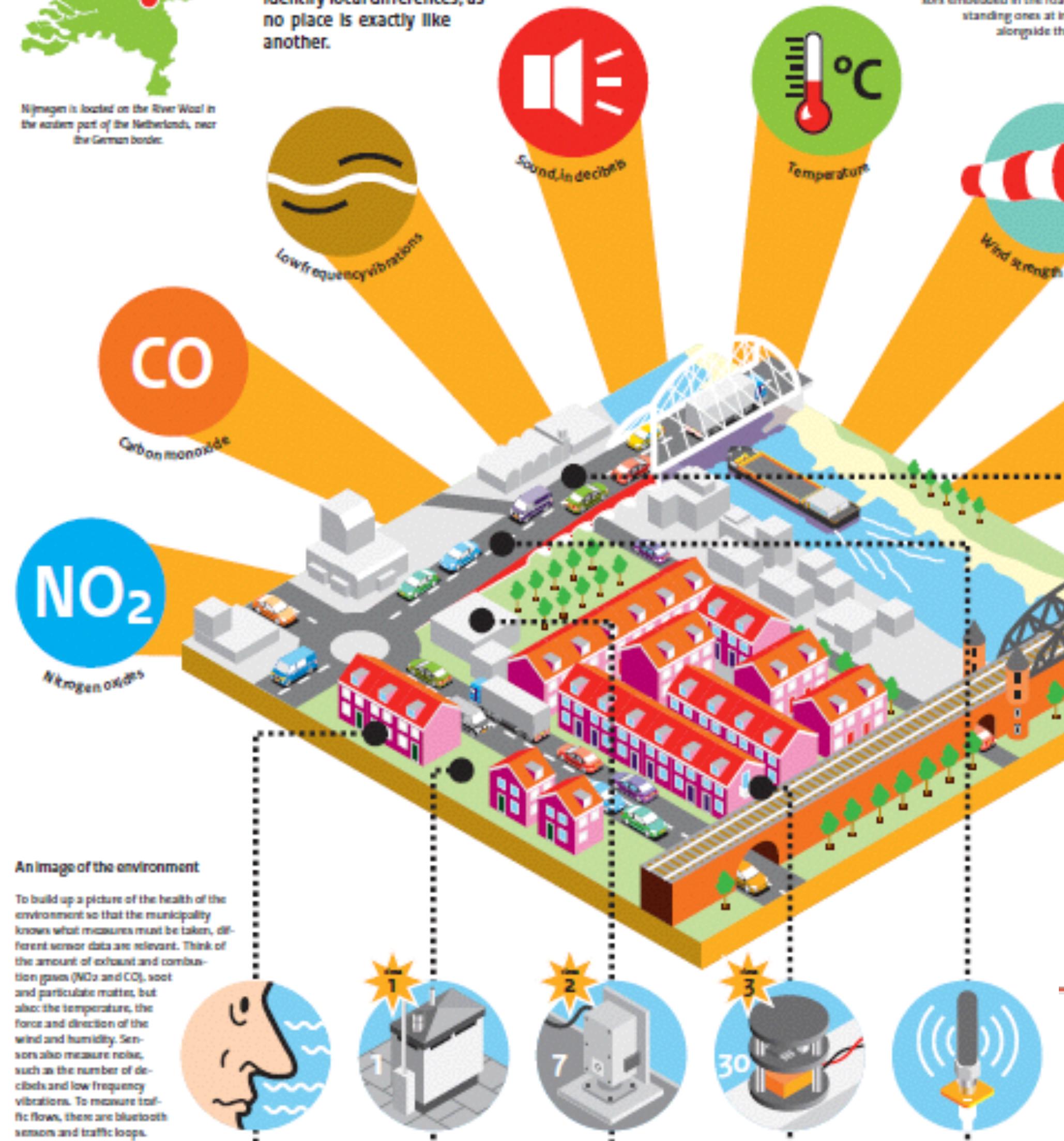


Case: Environmental health in Nijmegen



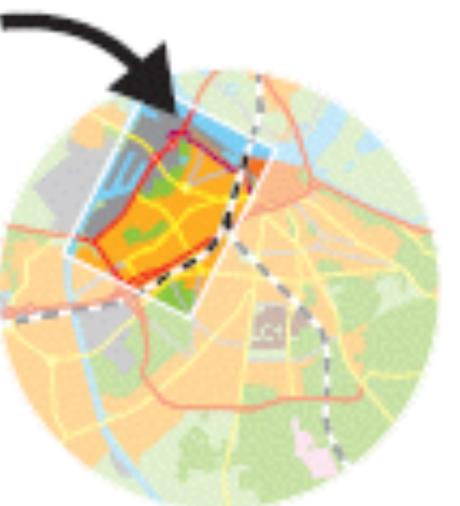
The environment is very important for people's health. That's why standards are set for the concentrations of pollutants. Sensors measure whether we keep within these standards using national monitoring networks. In addition, several municipalities and regions have their own sensors to identify local differences, as no place is exactly like another.

Nijmegen is located on the River Waal in the eastern part of the Netherlands, near the German border.



Nijmegen also monitors local environmental quality. With the advent of a new bridge and the construction of a ring road, the traffic situation in the western part of Nijmegen has changed. Developments in the port and the industrial area by the River Waal have been made, and residents in the nearby neighbourhood are worried about the health of their environment. The municipality is taking their concerns seriously and has placed sensors in the neighbourhood to measure the air quality and noise level.

Nijmegen also wants to collate reports about bad odours. To monitor traffic flows, the municipality uses sensors embedded in the road and free standing ones at intervals alongside the road.



matter sensor units and, as part of a research project, the Radboud university has distributed a 'warm' of thirty simple sensors among residents. One of the questions posed by this research is

From Pilot to Platform

Radboud Universiteit



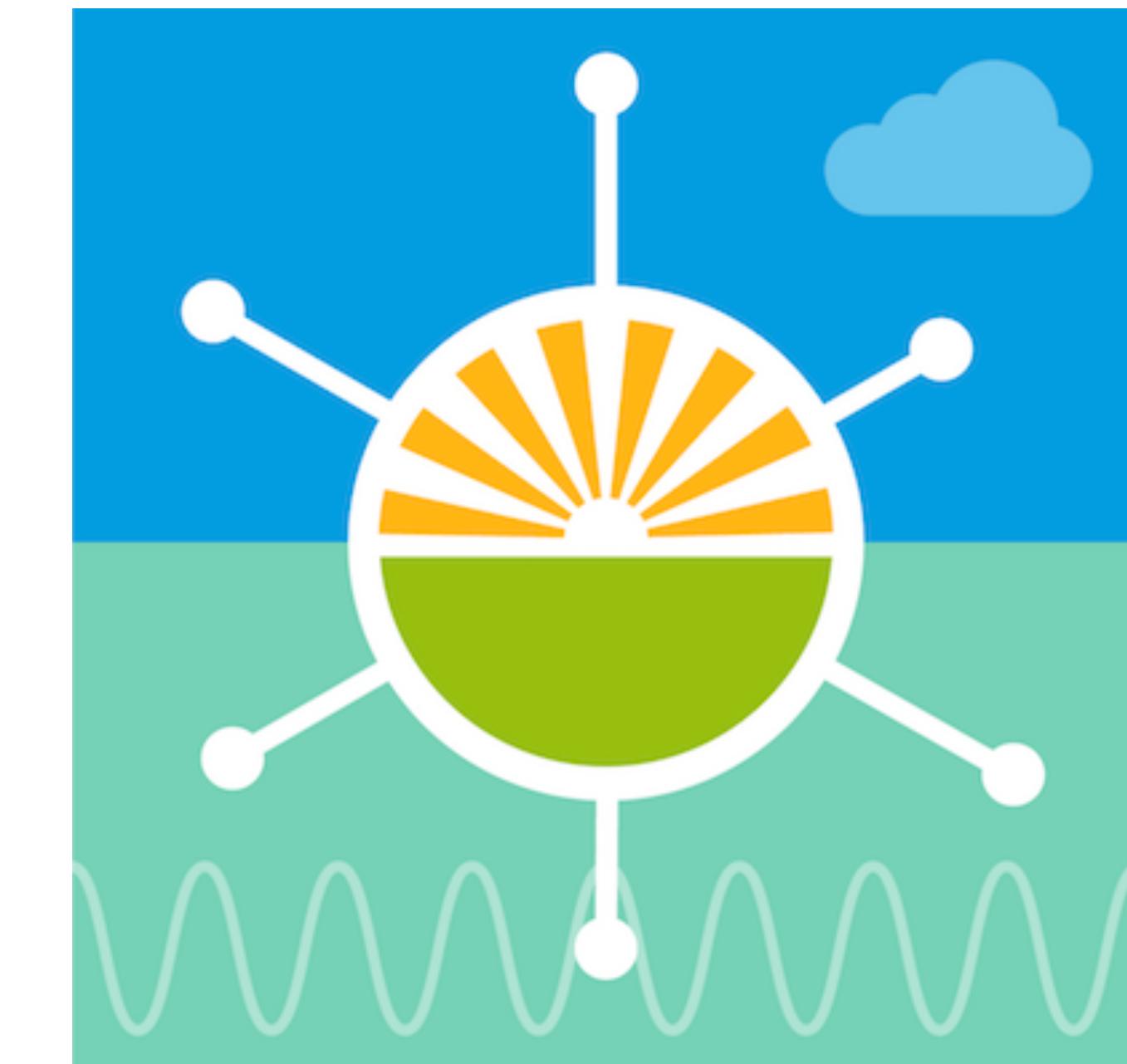
Intemo

INNOVATIVE TECHNOLOGY IN MOTION

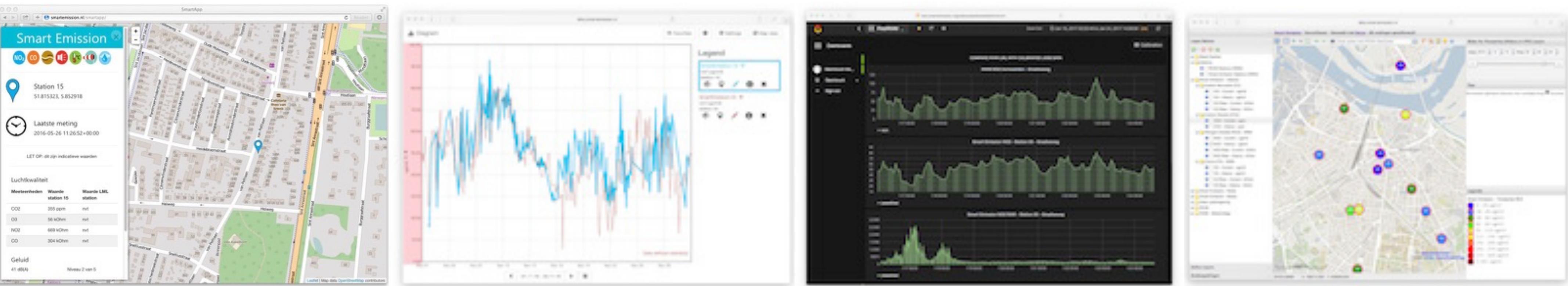


Rijksinstituut voor Volksgezondheid
en Milieu
Ministerie van Volksgezondheid,
Welzijn en Sport

powered by



Viewers



Stet!

Services

ETL

Sensors



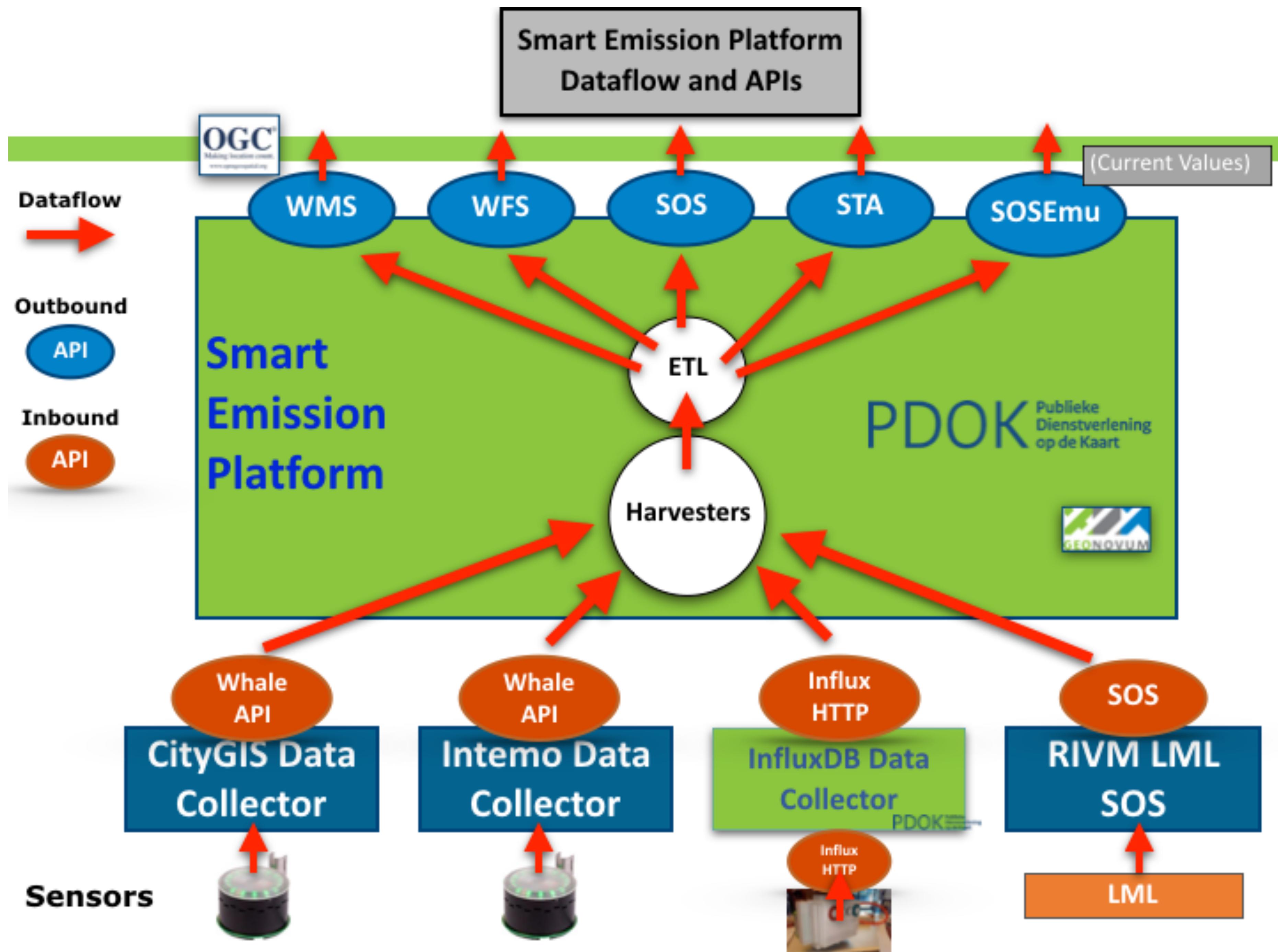
Smart Emission Platform

Smart Emission Platform Summary

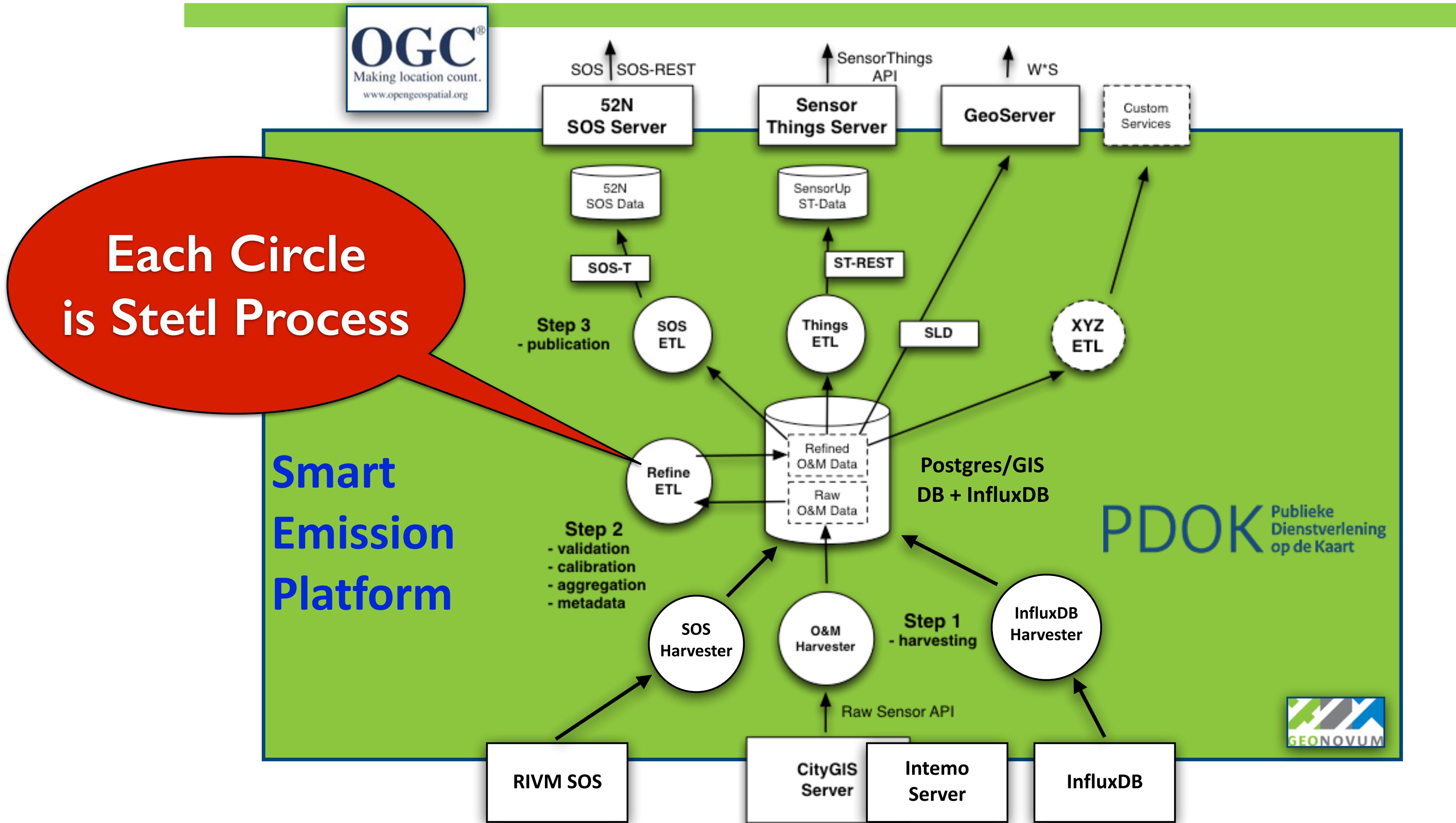
- Air Quality + Noise + Meteo Sensors hosted by citizens
- Raw sensor data Harvesting to SE Platform
- Validation, Calibration, Aggregation
- Publication to Web Services (WMS, WFS, SOS, SensorThings API)
- data.smartemission.nl

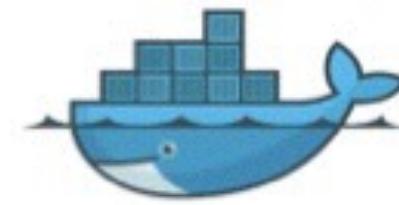
All ETL

(Harvesting, Validation, Calibration, Aggregation, Publication)
done with Stetl

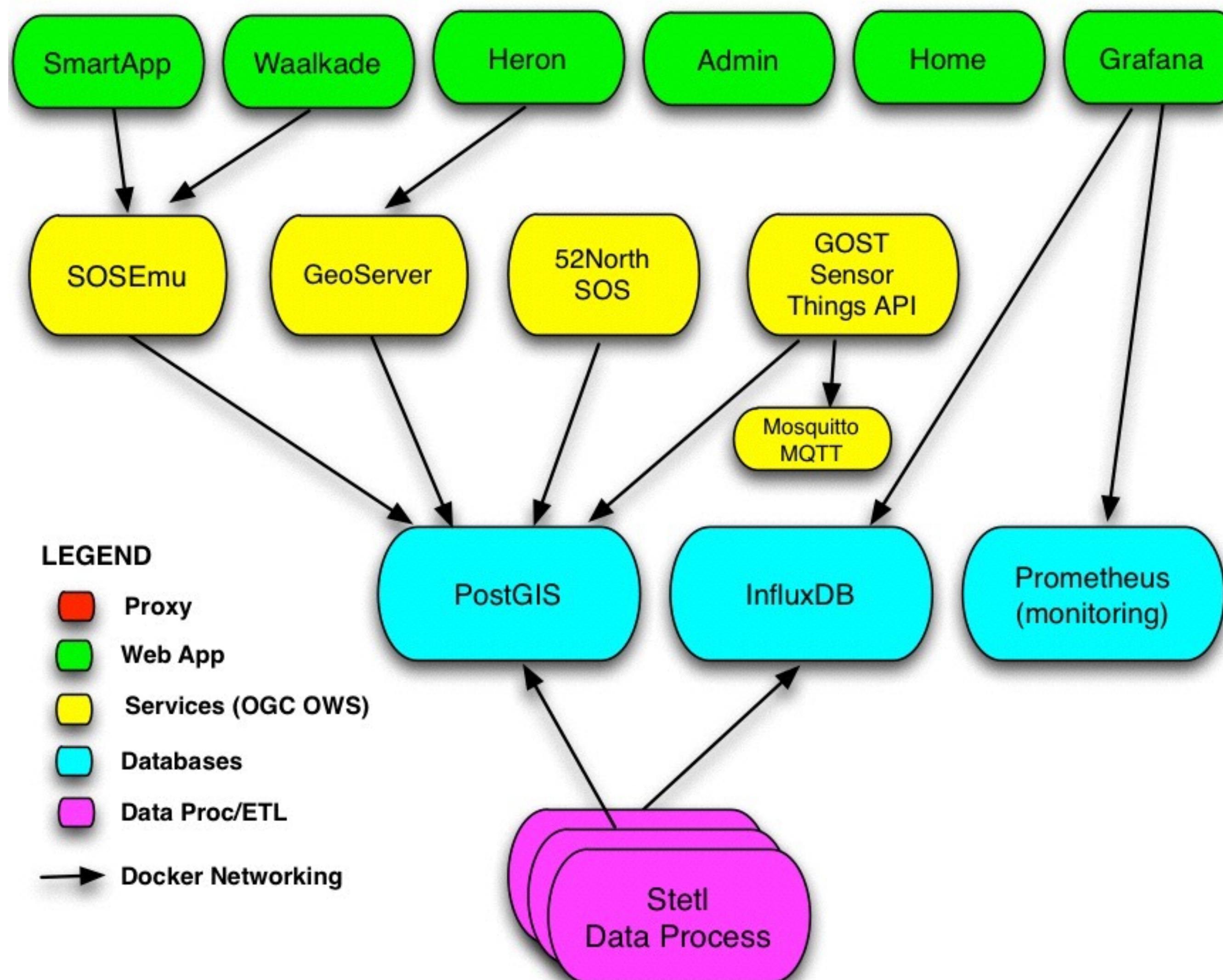


Data Architecture with multi-Step ETL





Smart Emission Docker Deployment



Smartem ETL Refiner

Project

1: Project

2: Structure

extractor.cfg x sospublisher.cfg x refiner.cfg x Dockerfile x

```

1  # Smart Emission Data Refiner ETL - Stel config
2  #
3  # Just van den Broecke - 2016-2018
4  #
5  # This config reads the raw timeseries measurements from the DB, refines these
6  # and writes results in the measurements DB.
7  #
8  # The main Stel ETL chain
9  [etl]
10 chains = input_raw_sensor_db | refine_filter| (output_postgres_insert) (component_sieve|output_influxdb_write)
11 # chains = input_raw_sensor_db/refine_filter/output_postgres_insert
12
13 # read raw data from timeseries table
14 [input_raw_sensor_db]
15 class = smartem.rawdbinput.RawDbInput
16 host = {pg_host}
17 database = {pg_database}
18 user = {pg_user}
19 password = {pg_password}
20 schema = {pg_schema_raw}
21 table = timeseries
22 output_format = record
23 last_gid_query = SELECT gid_raw from smartem.refined.refiner_progress
24 max_input_records = {refiner_max_input_records}
25 gids_query = SELECT gid from timeseries WHERE gid > ?d AND gid <= ?d AND complete = TRUE ORDER BY gid
26 data_query = SELECT * from timeseries WHERE gid = ?d
27 read_once = {refiner_raw_read_once}
28
29 # Refines raw records for specified sensor names
30 [refine_filter]
31 class = smartem.refiner.refinefilter.RefineFilter
32 sensor_names = temperature,humidity,pressure,noiseavg,noisalevelavg,co2,o3,co,no2,o3raw,coraw,noraw,no2raw,pml0,pm2_5
33 host = {pg_host}
34 database = {pg_database}
35 user = {pg_user}
36 password = {pg_password}
37 schema = {pg_schema_calibrated}
38 process_name = refiner
39
40 # Sieve out only gas-components
41 [component_sieve]
42 class = filters.sieve.AttrValueRecordSieve
43 input_format = record_array
44 output_format = record_array
45 attr_name = name
46 attr_values = o3,co,no2,co2
47
48 # for testing/debugging
49 [output_std]
50 class = outputs.standardoutput.StandardOutput
51
52 # Insert file records
53 [output_postgres_insert]
54 class = outputs.dboutput.PostgresInsertOutput
55 input_format = record_array
56 host = {pg_host}
57 database = {pg_database}
58 user = {pg_user}
59 password = {pg_password}
60 schema = {pg_schema_refined}
61 table = timeseries
62 replace=True
63
64 # Write records to InfluxDB server
65 [output_influxdb_write]
66 class = smartem.influxdboutput.InfluxDbOutput
67 input_format = record_array
68 method = POST

```

Run Smartem ETL Refiner

JavaEE App

favorites

Spe

fact

log

GUI ;-)

I. Dockerfile

```
1 FROM geopython/stetl:1.2
2
3 # Use the standard Stetl Docker Image with some additional packages
4 LABEL maintainer="Just van den Broecke <justb4@gmail.com>"
5
6 ENV ADD_PYTHON_DEB_PACKAGES="python-dev python-scipy python-seaborn python-matplotlib" \
7     ADD_PYTHON_PIP_PACKAGES="wheel Geohash influxdb scikit-learn==0.18"
8
9 # Not now see: https://github.com/pypa/pip/issues/5240
10 # RUN pip install --upgrade pip
11 RUN apt-get update && apt-get --no-install-recommends install -y \
12     ${ADD_PYTHON_DEB_PACKAGES} \
13
14     && pip install ${ADD_PYTHON_PIP_PACKAGES} \
15     && apt-get purge -y python-dev \
16     && apt autoremove -y \
17     && rm -rf /var/lib/apt/lists/*
18
19 # Copy relevant files to work dir
20 ADD config /work/config
21 ADD smartem /work/smartem
22 ADD options/example.args /work/options/default.args
23 ADD scripts /work/scripts
24
25 WORKDIR /work
26
27 ENTRYPOINT ["/work/scripts/entry.sh"]
28
29
```

2. entry.sh

```
1#!/bin/bash
2#
3# Generic ETL entry: override with specific ETL config and args file.
4#
5WORK_DIR=/work
6cd ${WORK_DIR}
7
8# Shorthand
9function log() {
10    echo "entry.sh: $1"
11}
12
13# Shorthand
14function error() {
15    log "$@"
16    exit -1
17}
18
19# ETL Process can be specified via env var or argument
20if [ "${ETL_PROCESS}"x = "x" ];
21then
22    ETL_PROCESS="$1"
23fi
24
25# ETL Process can be specified via env var or argument
26if [ "${ETL_PROCESS}"x = "x" ];
27then
28    error "Error: no ETL process specified"
29fi
30
31ETL_CONFIG_FILE=config/${ETL_PROCESS}.cfg
32if [ ! -f ${ETL_CONFIG_FILE} ];
33then
34    error "Error: cannot find ETL config file: ${ETL_CONFIG_FILE}"
35else
36    log "OK using config file: ${ETL_CONFIG_FILE}..."
37fi
38
39# Start Stetl with config file
40export PYTHONPATH=${WORK_DIR}:${PYTHONPATH}
41stetl -c ${ETL_CONFIG_FILE} -a options/default.args
42
43
44
```

3. Run ETL Process

```
$ docker run -t smartemission/se-stetl:1.0.13 refiner
```

Thank You !

www.stetl.org

github.com/geopython/stetl