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

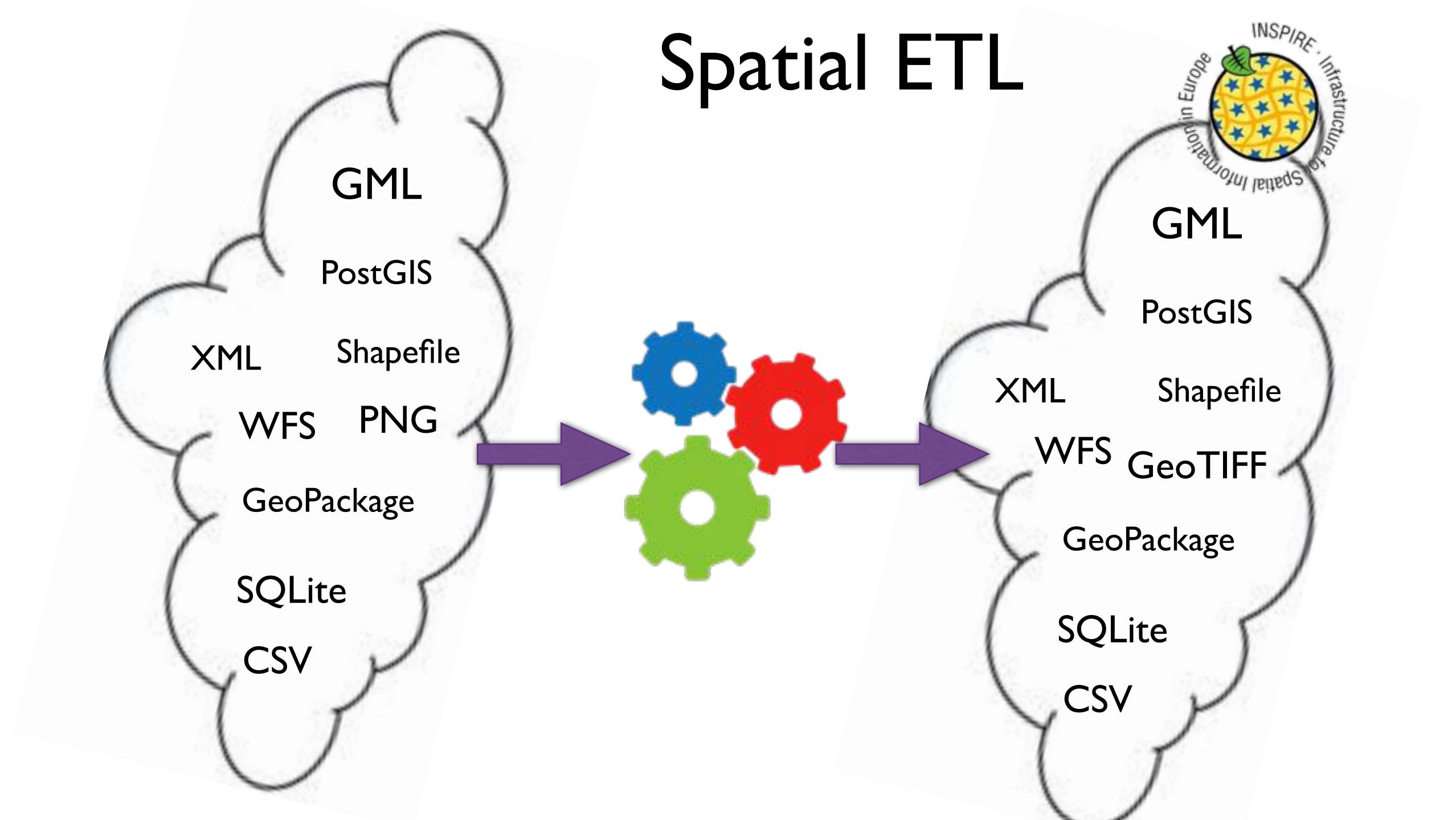


Agenda

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

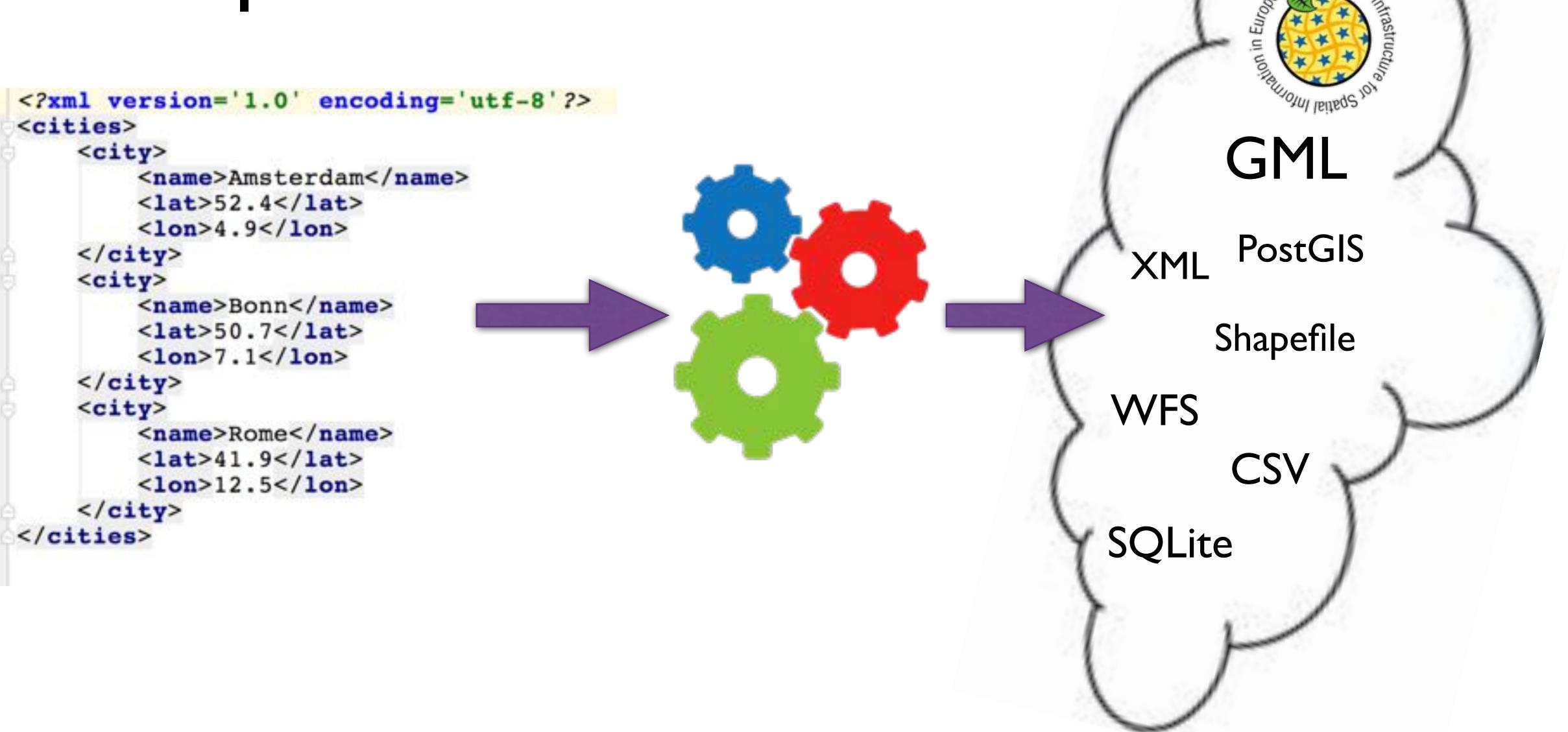
Spatial ETL

ETL - Extract Transform Load 一至三



Spatial ETL

Example non-standard source



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, Daylon 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, SOI 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







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



20

UNITE 21









12:49 PM - 25 Jan 2016









Plenty of Tools...







ogr2ogr

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

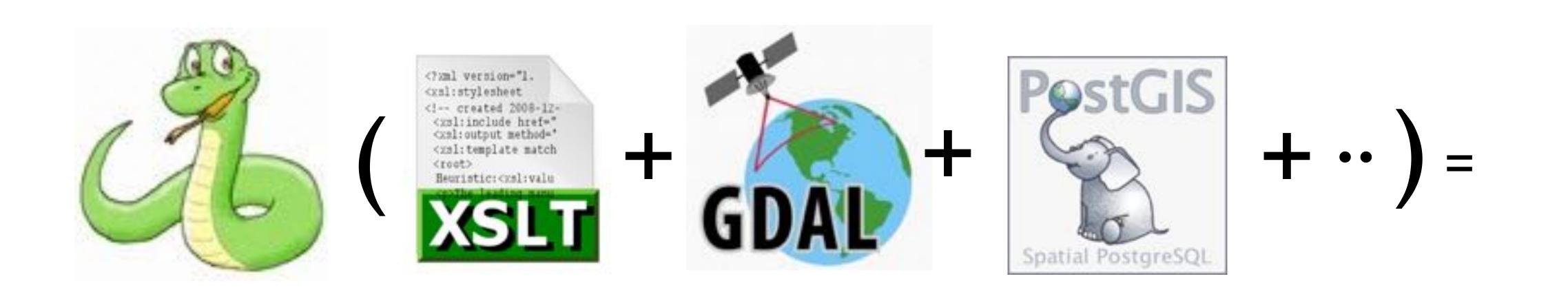
FOSS ETL - How to Combine Components?







Solution: Add Python to the Equation



Stetl

Stetl

Simple Streaming Spatial Speedy ETL



And what about?



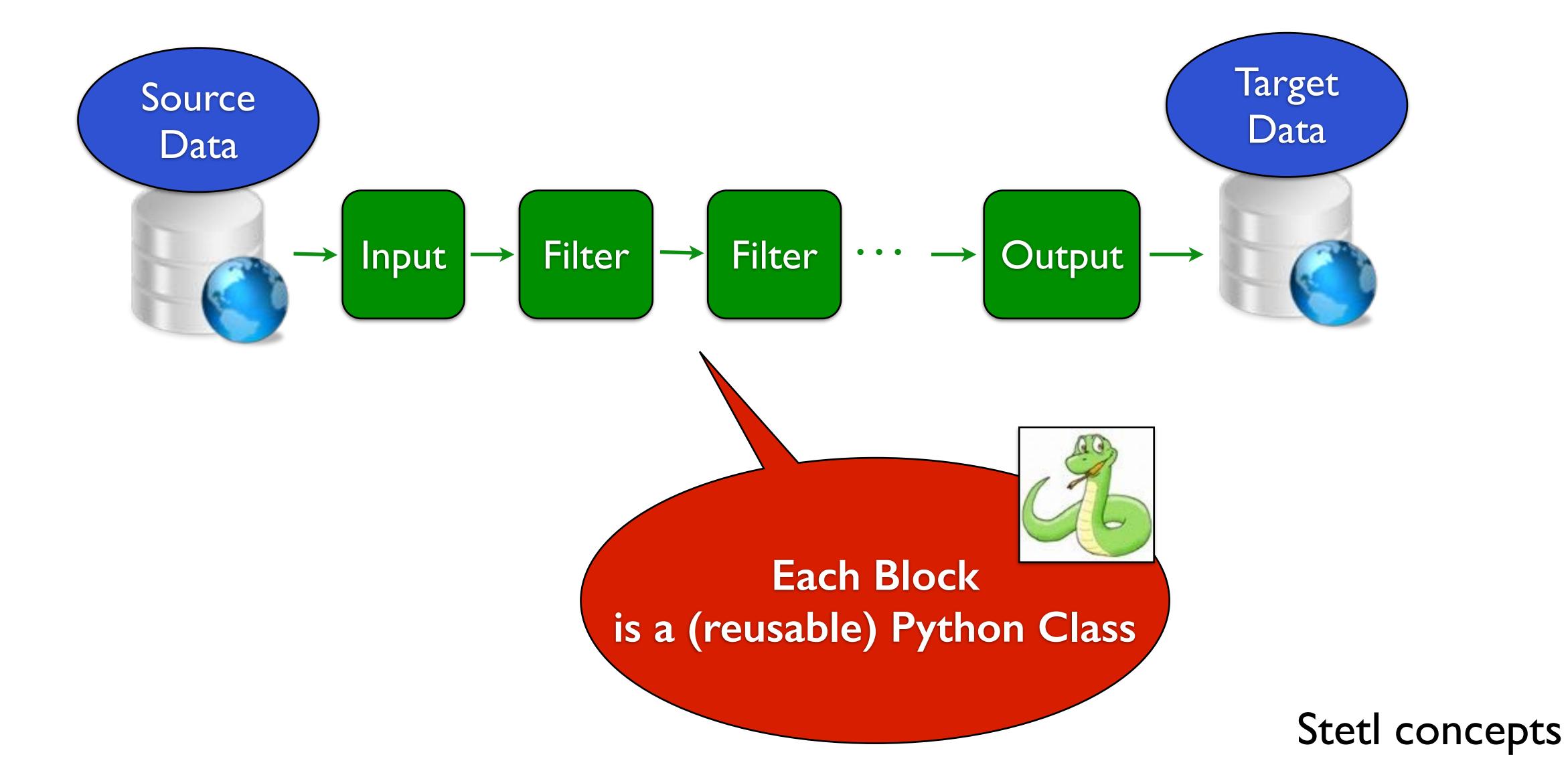




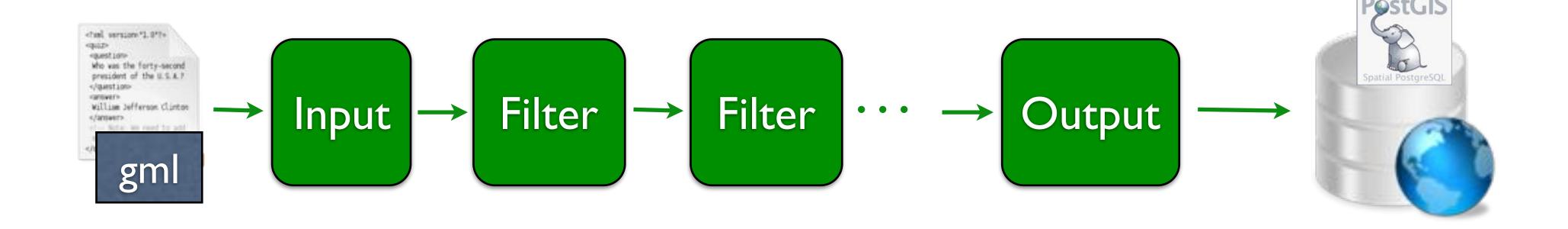


Stetl Concepts

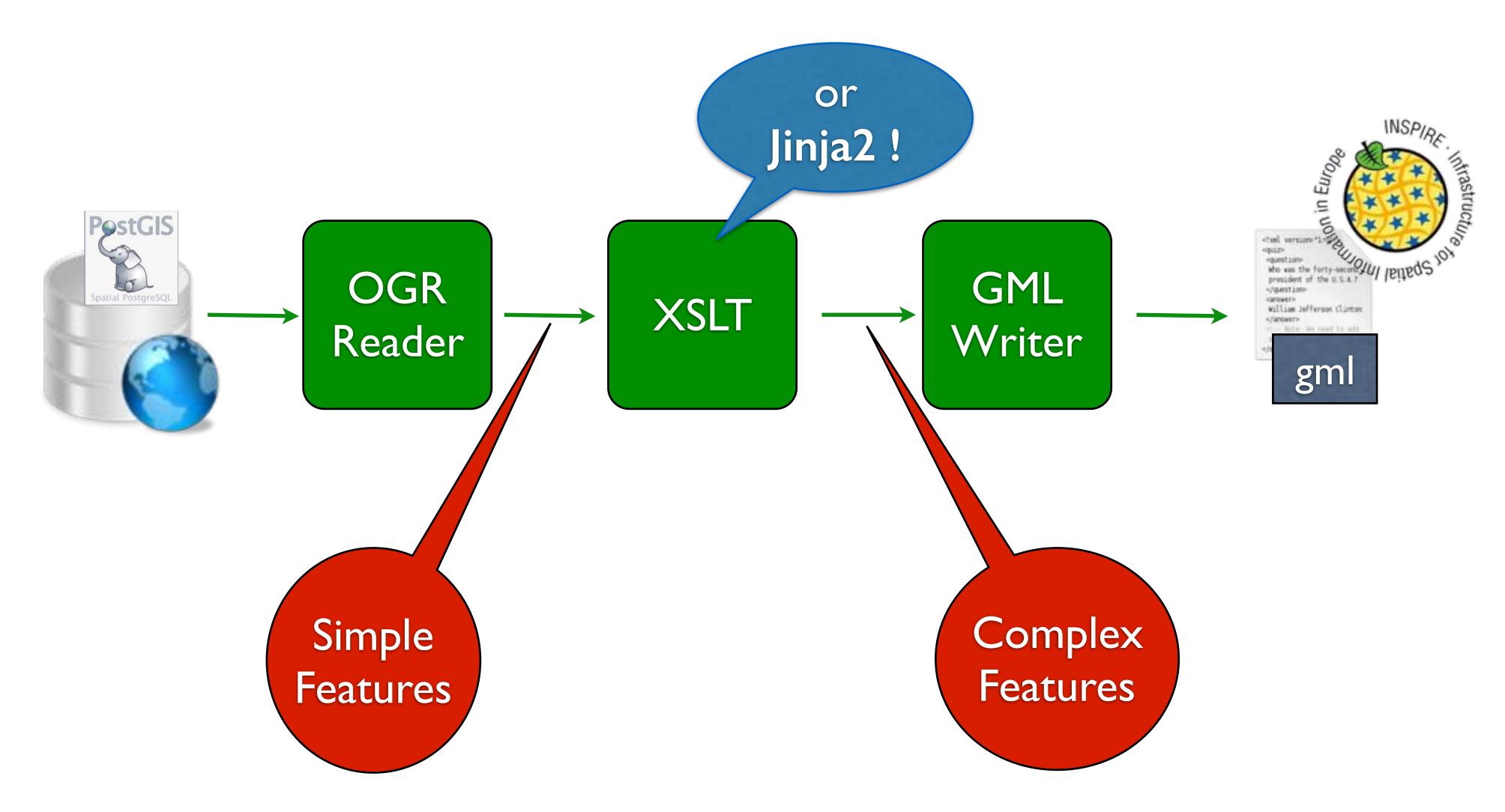
Process Chain



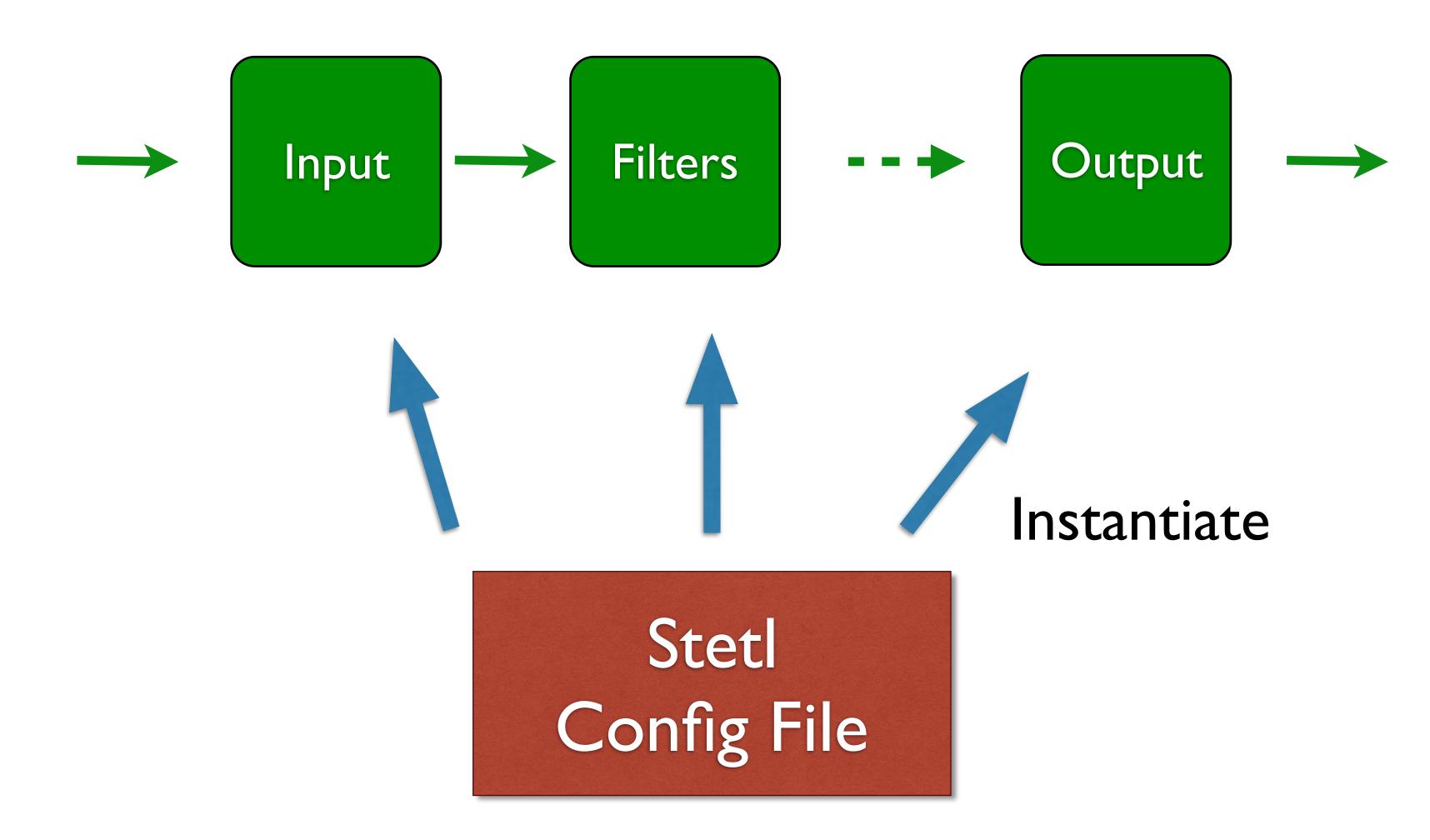
Process Chain - GML to PostGIS



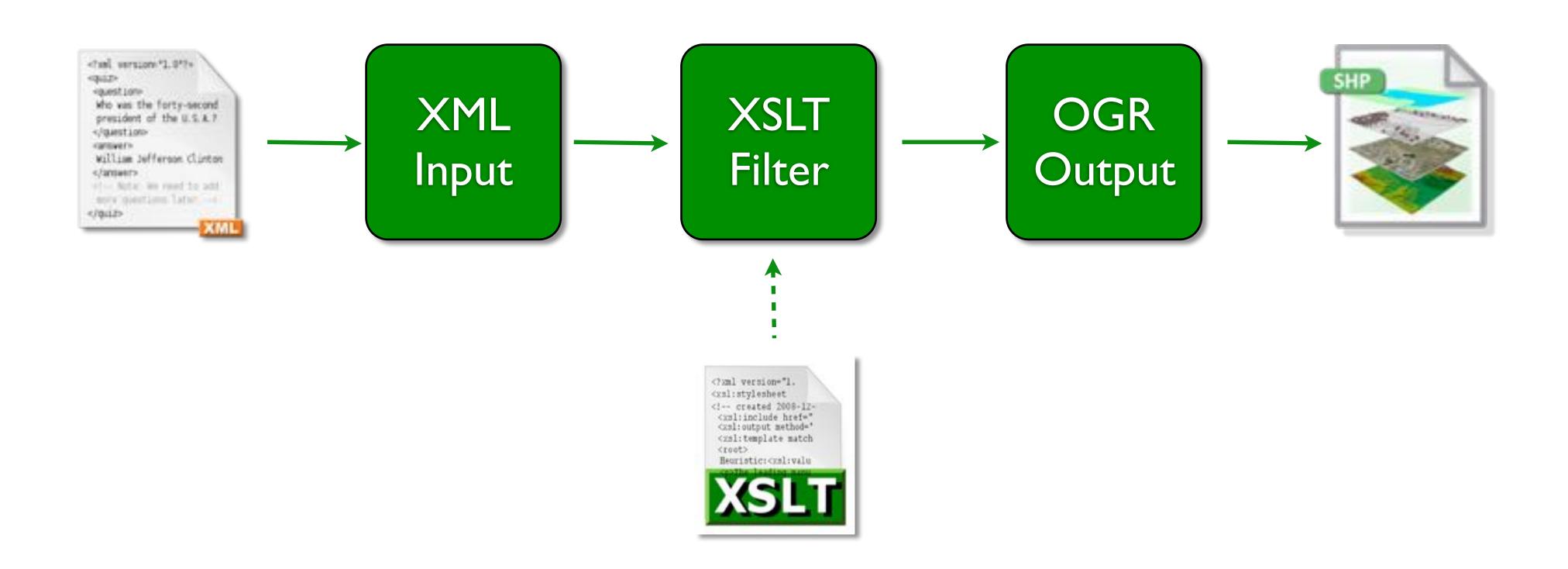
Example: Data Model Transform



Process Chain - How?

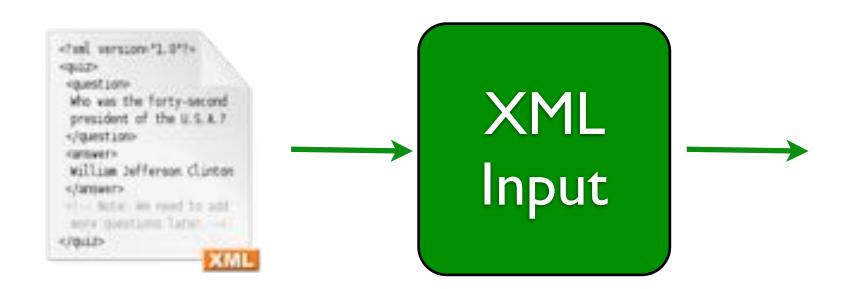


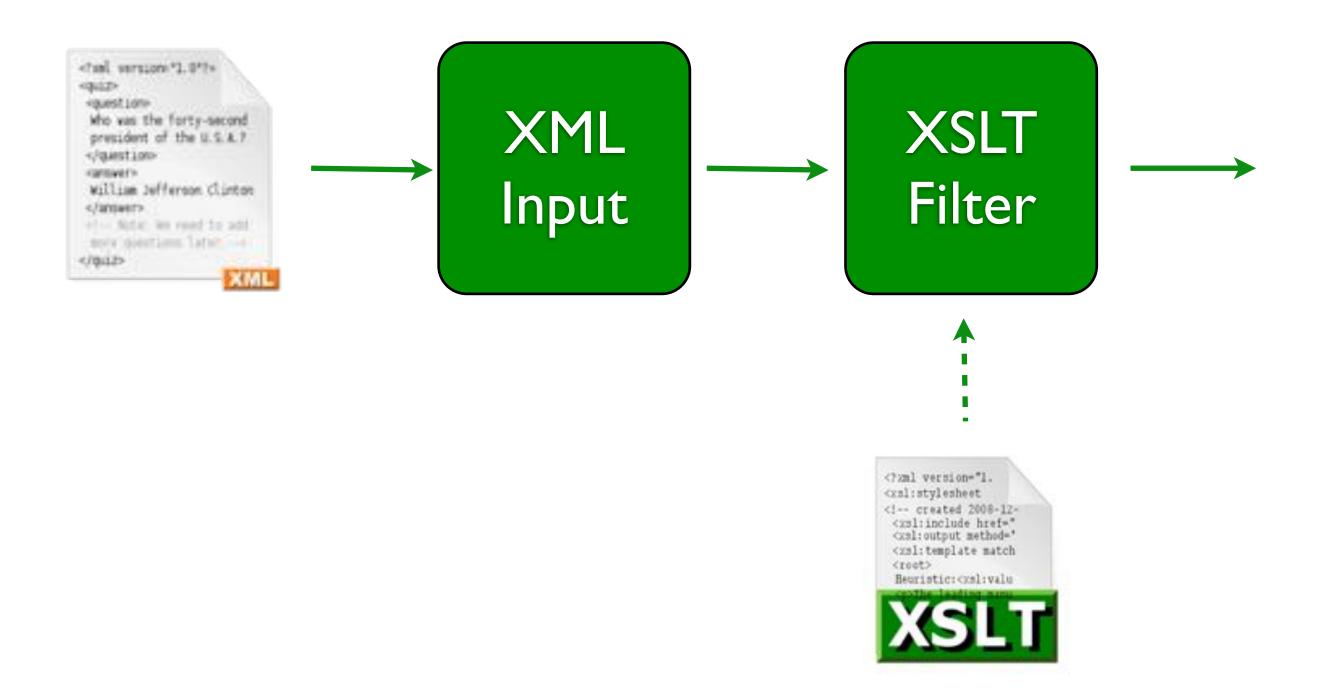
Example: XML File to Shapefile



```
<
<cities>
      <city>
            <name>Ansterdan</name>
            <lat>52.4</lat>
            <los>4.9</los>
      </city>
      <pity>
            <mame>Boom</pame>
            <lat>50.7</lat>
            <los>7.1</los>
      </city>
      <city>
            <name>Rome
            <lat>41.9</lat>
            <loa>12.5</loa>
      </city>
</cities>
```

The Source File

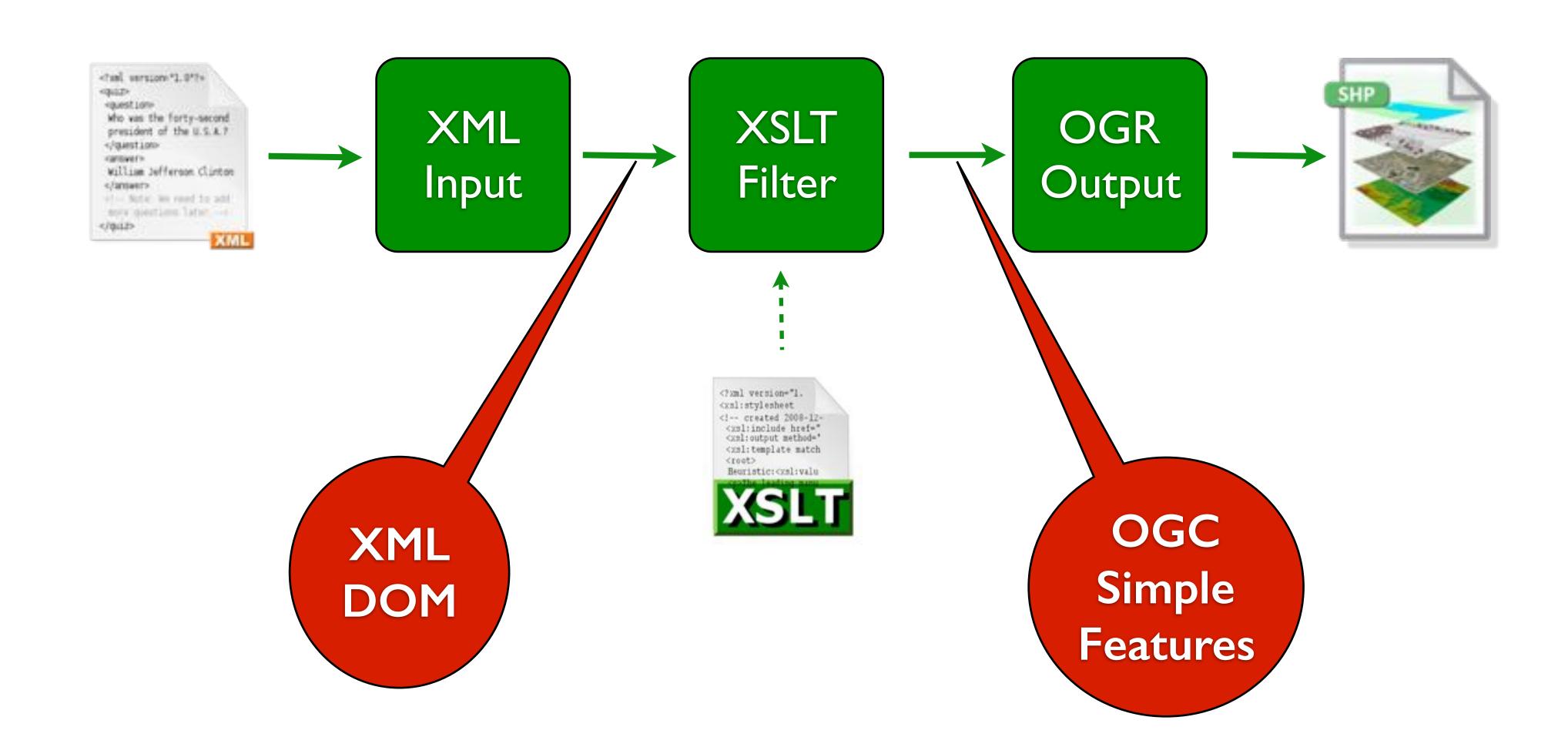




```
<xsl:template natch="/">
       -cogr:FeatureCollection
             xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
             xmlns:ogr="http://ogr.maptools.org/"
             xnlns:gnl="http://www.apengis.net/gnl"
             xsi:schemalocation="http://ogr.maptools.org/ ../gmlcities.xsd http://www.opengis.net/gml http://schemas
          sgml:boundedBy>
            <gnl; Box>
             <gml:coord><gml:X>-188.0
             <gnl:coord><gnl:X>180.0/gnl:X><gnl:Y>90.0/gnl:Y>
            </gral:Box>
          </graft/boundedBy>
           -i- Loop through all cities --
          <xsl:apply-templates/>
       </xsl:template>
   <!- Make each city an ografestureMember. ->
                                                      Prepare XSLT Script
   -xxsl:template match="city">
       <gnl:featureMember>
          cogr:City>
              cogr|name>
                 -cxsl:value-of select="name"/>
             </ogr:name>
              <corr:geometry>
                 <gml:Point srsName="urn:ogc:def:crs:EP56:4326">
                    *gml:coordinates=*xsl:value-of select="lat"/>, <xsl:value-of select="lon"/=/gml:coordinates=</pre>
                  </gml:Point>
              </or>
          </ogr/City>
       </xsl:template>
</xsl:stylesheet>
```

```
Open recoordings assessed 4.4.4.90/pml reconstituates
      STARREST PRESENTS
    4/ogrigeometry>
                                          XSLT GML Output
      right (Potat scatterer urbrogg) det correstrate distinct
       riged incoordinates wenter. 7.7.14/gelicoordinates/
      47 god : Polintin
    TARGET GROWNINGS
  Cognicking.
```

Crombine to a restriction of the



XSLT

```
Process
          # Transform input xml to valid GML file and them to Shape.
          [etl]
          chains = input xml file | transformer xslt | output ogr shape
          [input_xml_file]
          class = inputs.fileinput.XmlFileInput
                                                        XML
          file path = input/cities.xml
                                                        Input
          [transformer_xslt]
Filter
          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 \
                                                 ogr2ogr
              -overwrite
              -f "ESRI Shapefile" \
                                                 Output
              -a srs epsg:4326
             output/gmlcities.shp \
             temp/gmlcities.gml
```

Chain

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
```

•

Stetl concepts

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

```
[et1]
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
```

Installing Stetl - PyPi

sudo pip install stetl

Deps

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

Installing Stetl - Docker



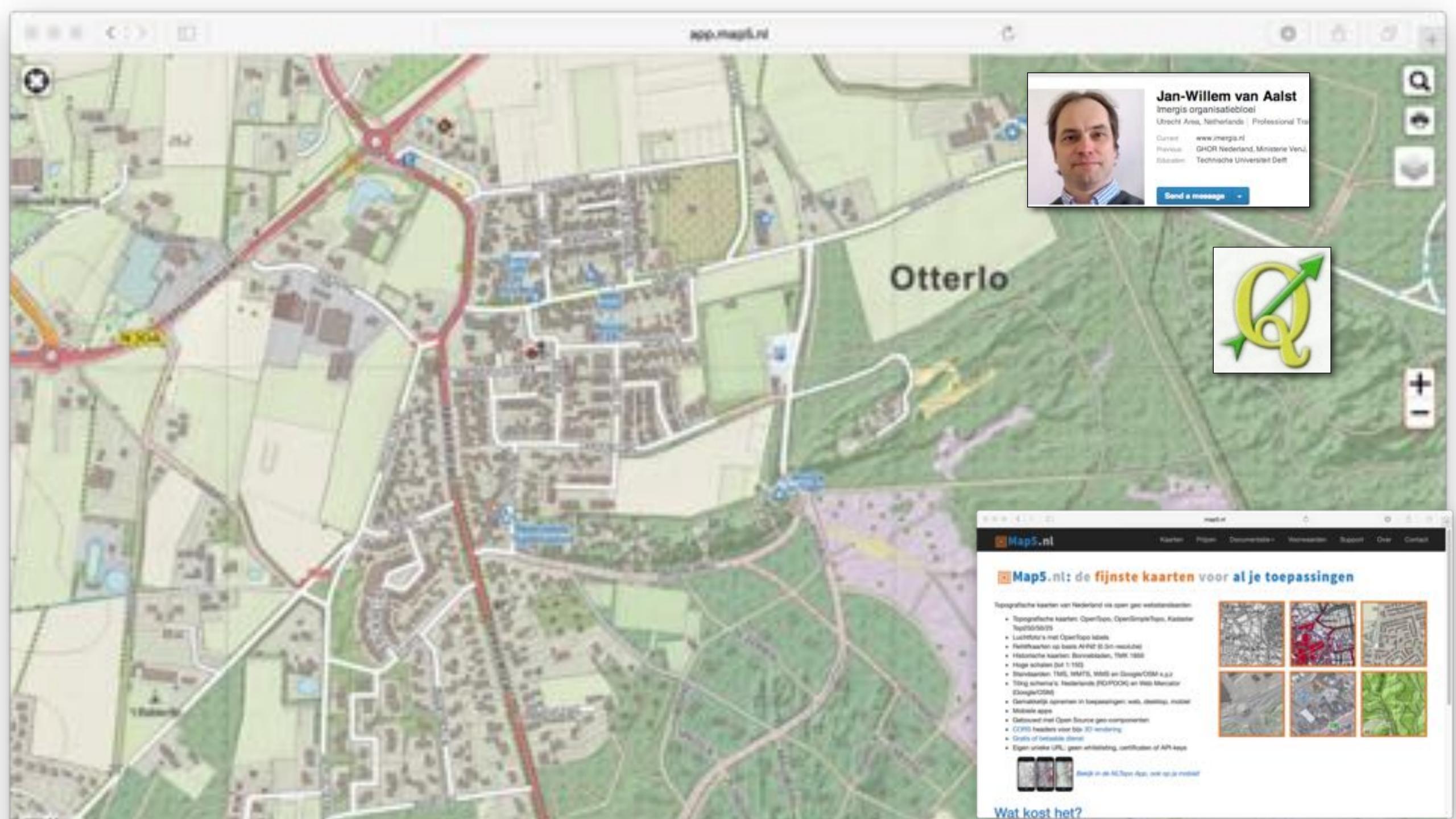
```
https://hub.docker.com/r/
geopython/stetl
```

Case I - NLExtract

What is NLExtract?

ETL for Dutch National Open Datasets http://nlextract.nl

- GML to PostGIS:
 BAG, BRT (Top I 0NL), BGT, BRK (Cadastral Parcels)
- Downloads <u>data.nlextract.nl</u>

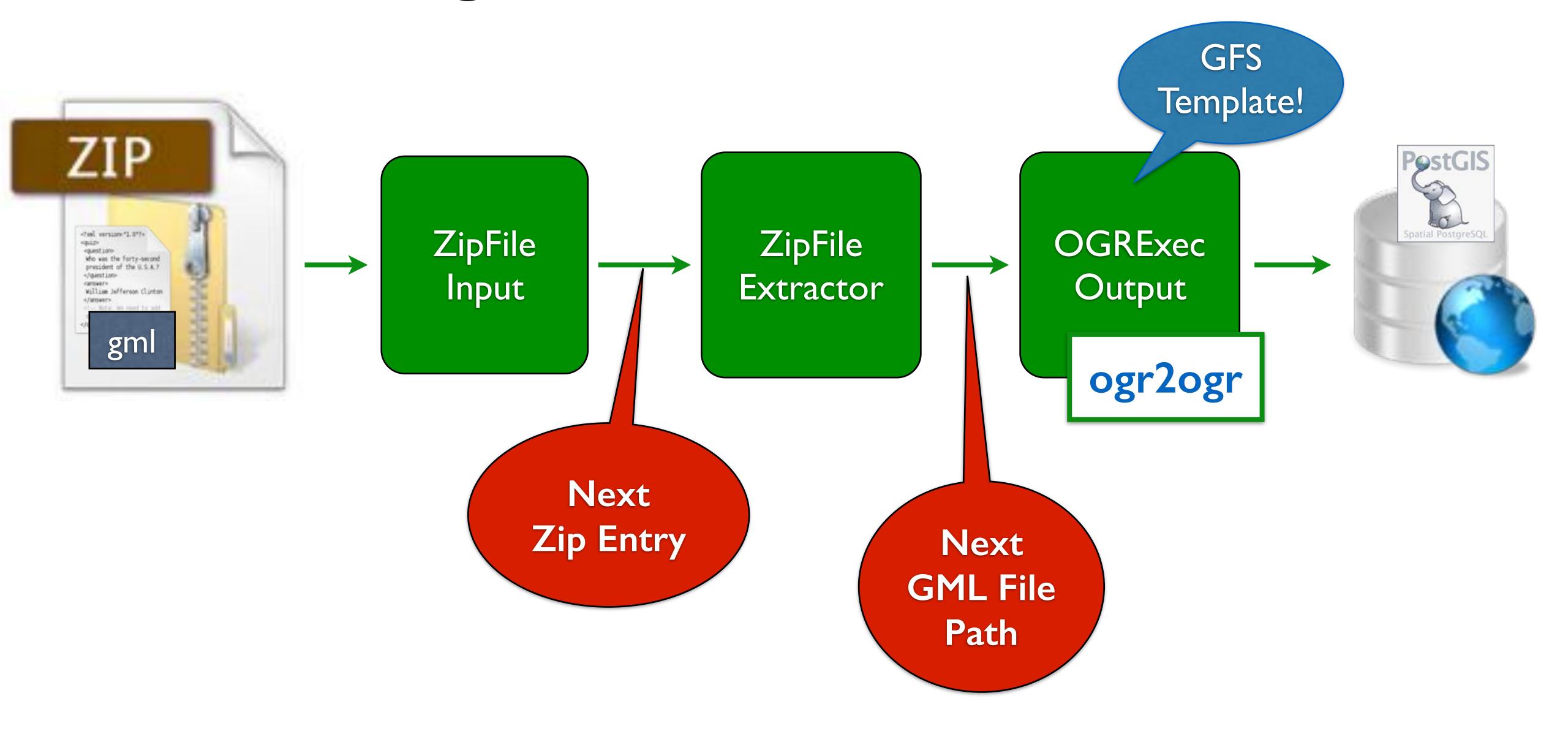


BRK - Input GML from PDOK

```
$ 1s -1h
total 3.6G
-rw-rw-r-- 1 mlx mlx 3.6G May 7 10:32 kadastralekaartv3-gml-nl-nohist.zip
$ unzip -1 kadastralekaartv3-gml-nl-nohist.zip
           Date Time
Length
                            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}
# qfs template
gfs_template = {gfs_template}
# miscellaneous ogr2ogr options
options = -append -gt 65536 {multi_opts} --config PG_USE_COPY NO --config CPL_ZIP_ENCODING CP437
```

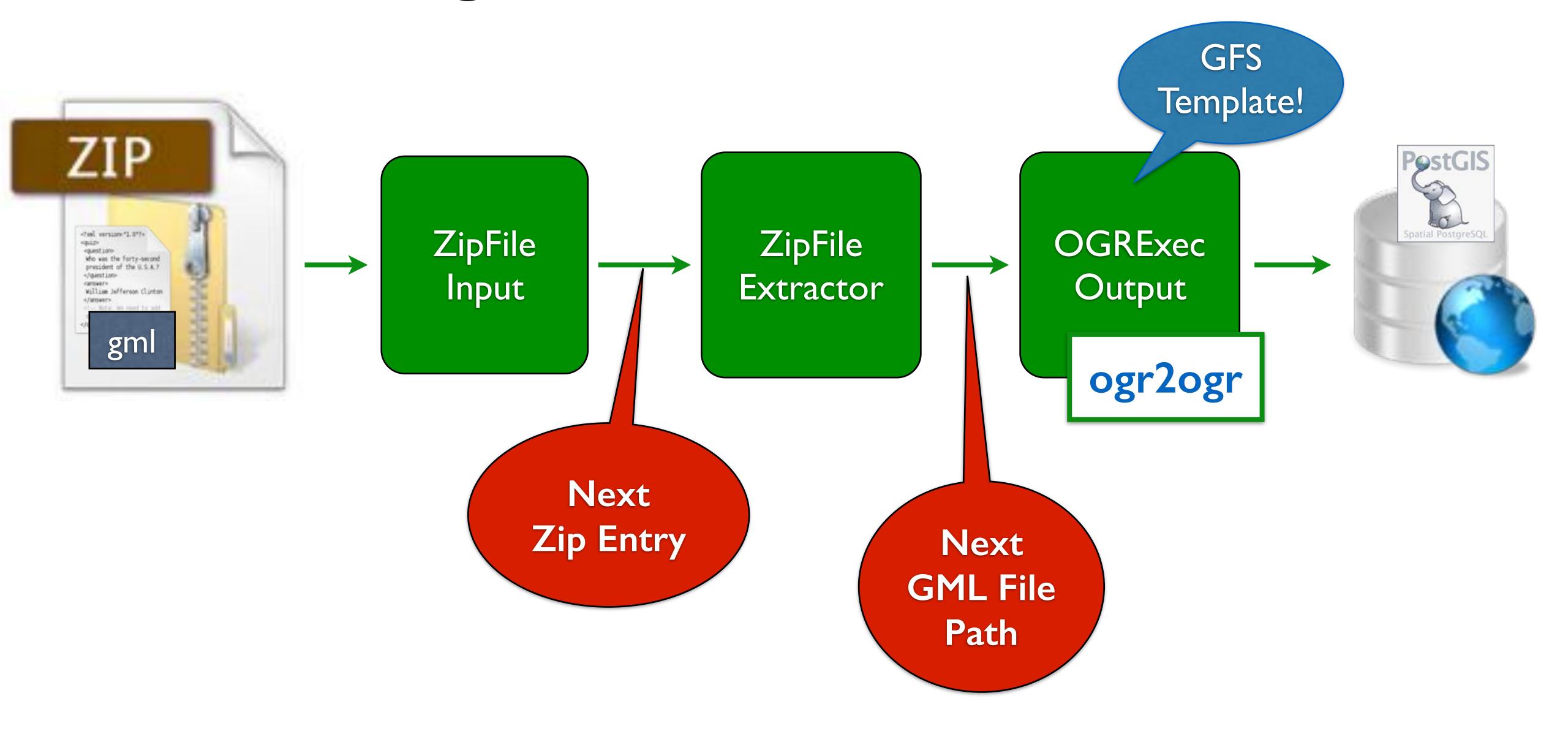
```
GFS for Table+Column Mapping
<GMLFeatureClassList>
    <GMLFeatureClass>
       <Name>Perceel</Name>
       <ElementPath>Perceel/ElementPath>
        <!--<GeometryType>1</GeometryType>-->
       <SRSName>urn:ogc:def:crs:EPSG::28992</SRSName>
       <PropertyDefn>
           <Name>gemeente</Name>
           <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

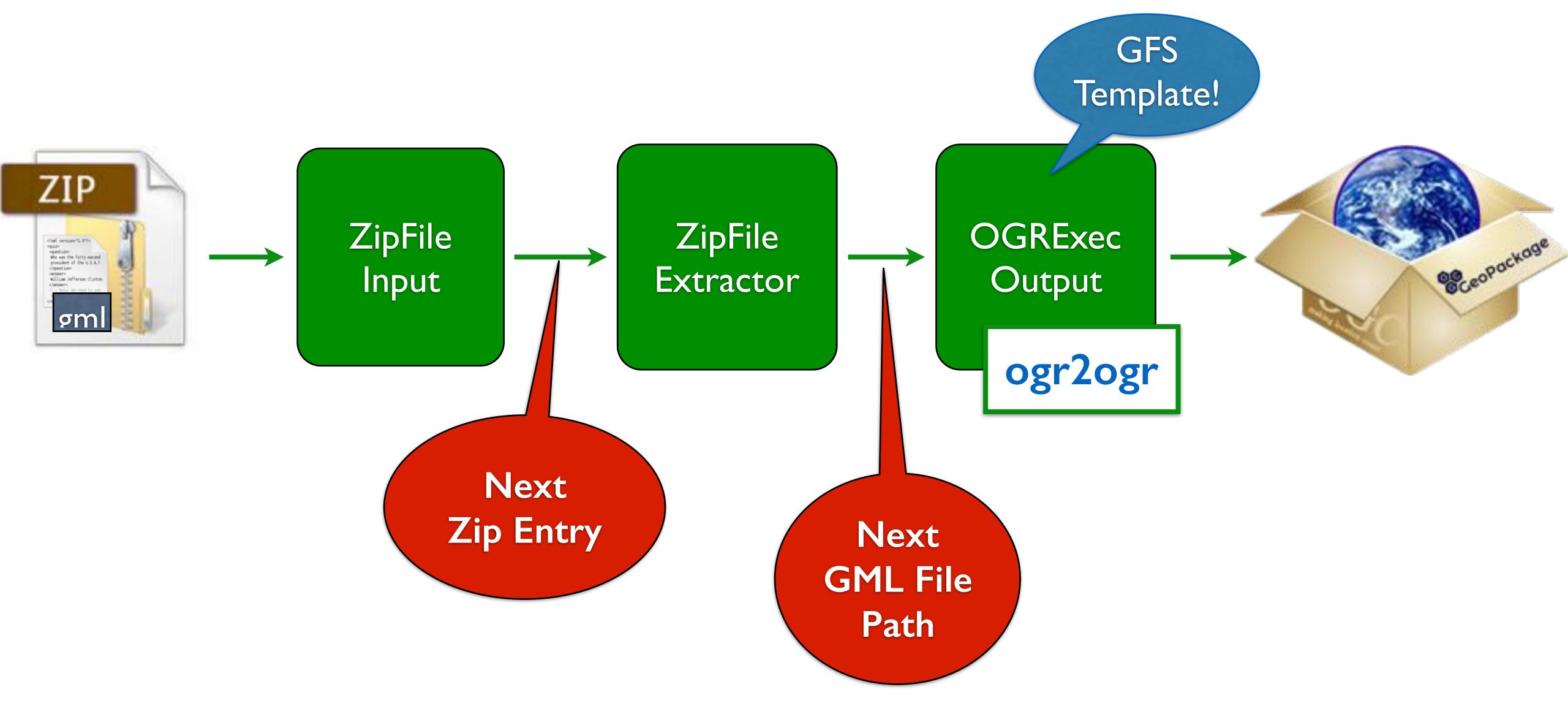
ogc_fid	gml_id	namespace	lokaalid	logischtijdstipontstaan	gemeente	sectie	perceelnummer	waarde
1	NL.IMKAD.KadastraalObject.120031707	NL.IMKAD.KadastraalObject	120031707	2009-03-27T23:59:59.000	AMR04	0	3322	178
2	NL.IMKAD.KadastraalObject.120118807	NL.IMKAD.KadastraalObject	120118807	2009-03-27T23:59:59.000	AMR04	Р	3907	446
3	NL.IMKAD.KadastraalObject.120077624	NL.IMKAD.KadastraalObject	120077624	2009-03-27T23:59:59.000	AMR04	0	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

perceelnummerrotatie	deltax	deltay	begrenzing
0	0	0	010300002040710000010000000C0000003D0AD7A3CF5E014
0	0	0	01030000204071000001000000060000006F1283C072D0014
0	0	0	010300002040710000010000000600000D34D6210DB6F014
-45.6	0	0	010300002040710000010000000900000046B6F3FD0C2D014
0	0	0	0103000020407100000100000000A000000C74B3789BE42014

Stetl Config BRK to PostGIS

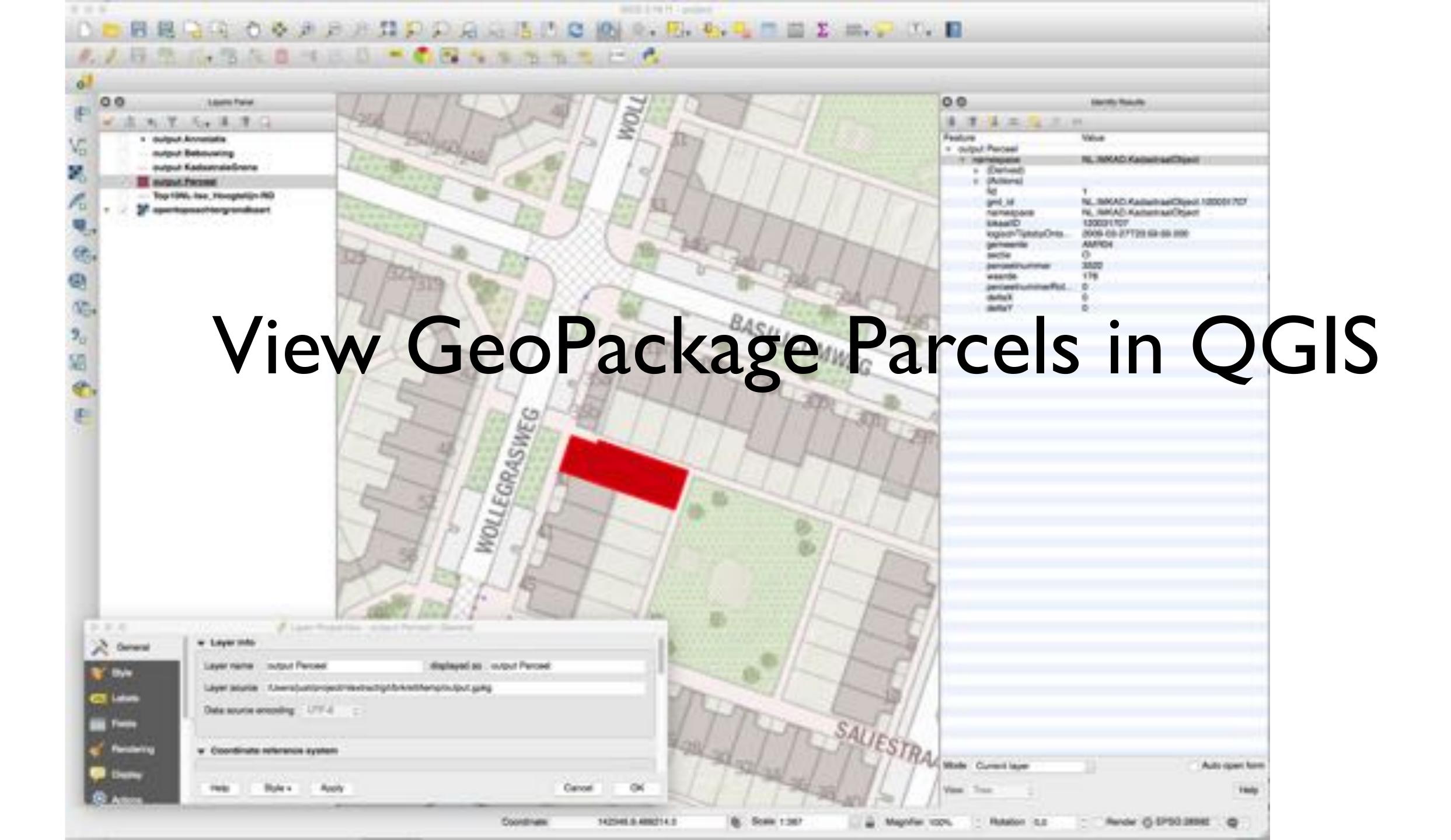


Stetl Config BRK to GeoPackage



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 datasedrce: 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
```

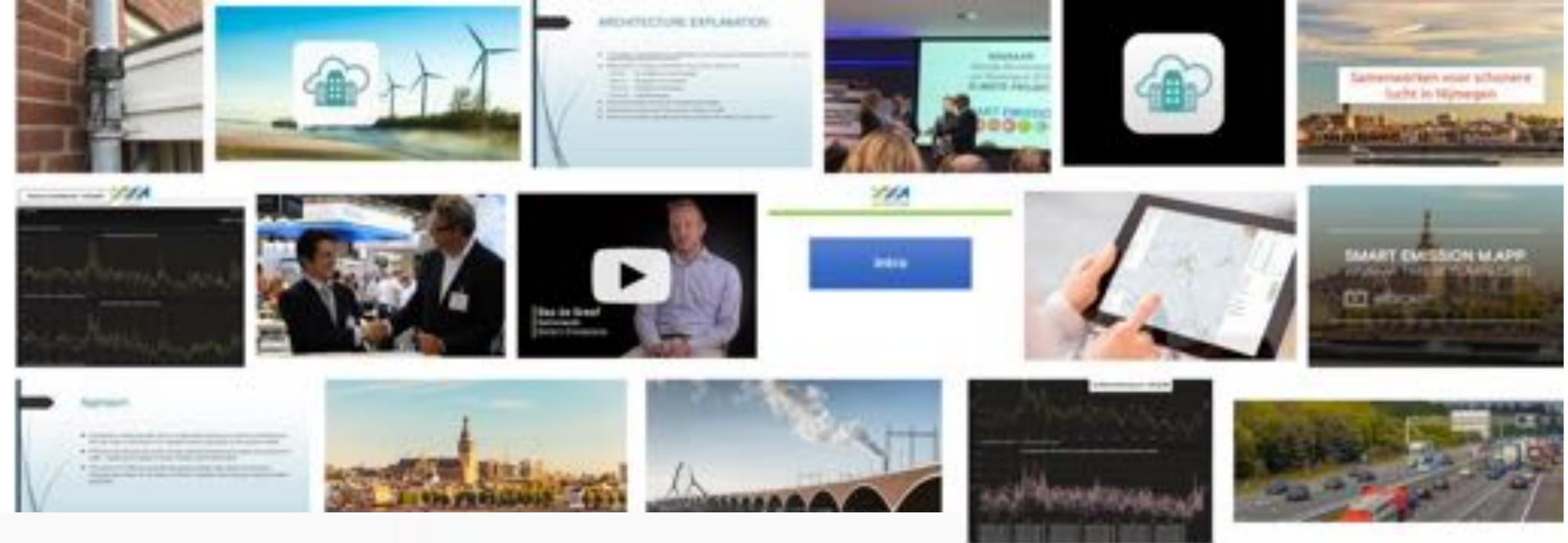


BRK .zip to both PostGIS + GeoPackage

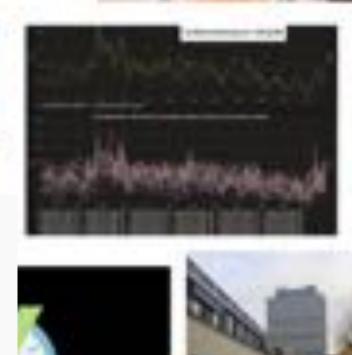
```
[etl]
chains = input_zip_file | extract_zip_file | (output_ogr2ogr_postgis) (output_ogr2ogr_gpkg)
# The source input ZIP-file(s) from dir, producing 'records' with ZIP file name and inner file names
[input zip file]
class=inputs.fileinput.ZipFileInput
file_path = {input_dir}
filename_pattern = {zip_files_pattern}
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 Stetl "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



From Pilot to Platform

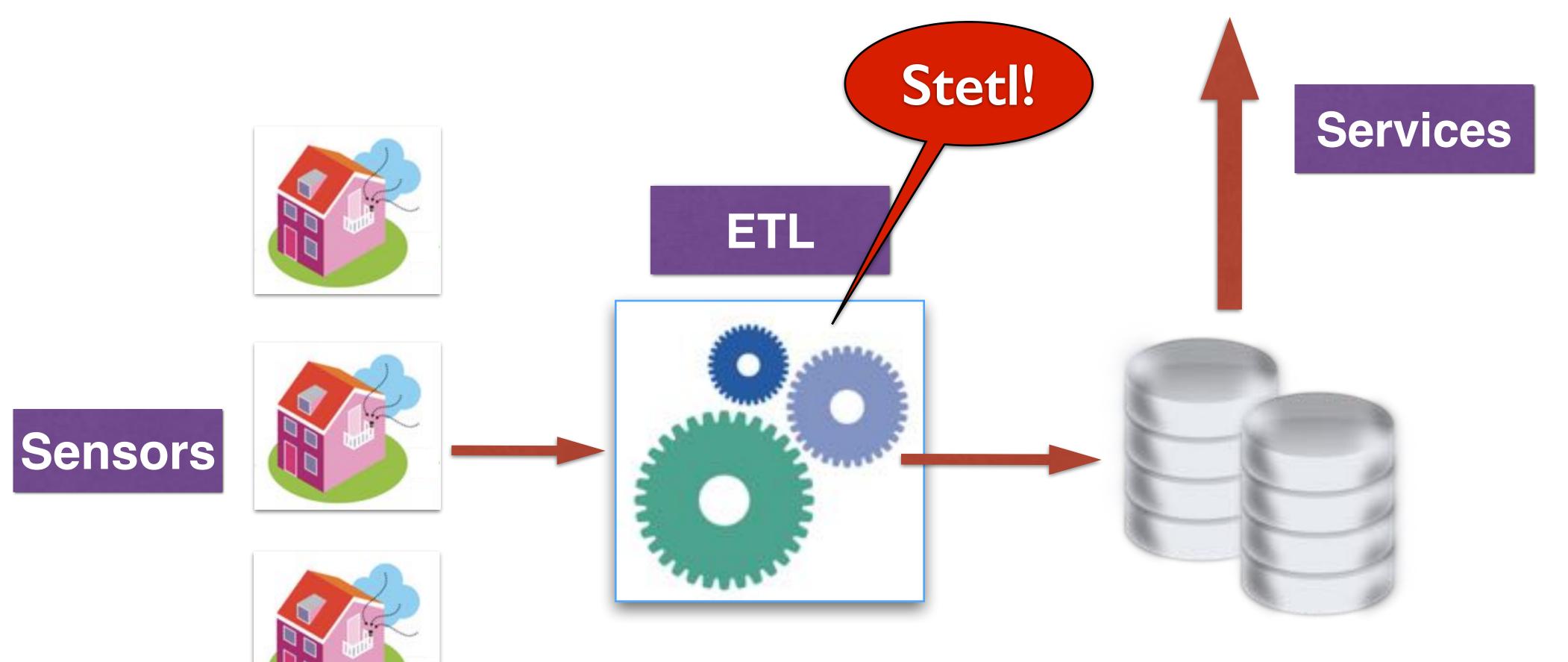






Viewers





Smart Emission Platform

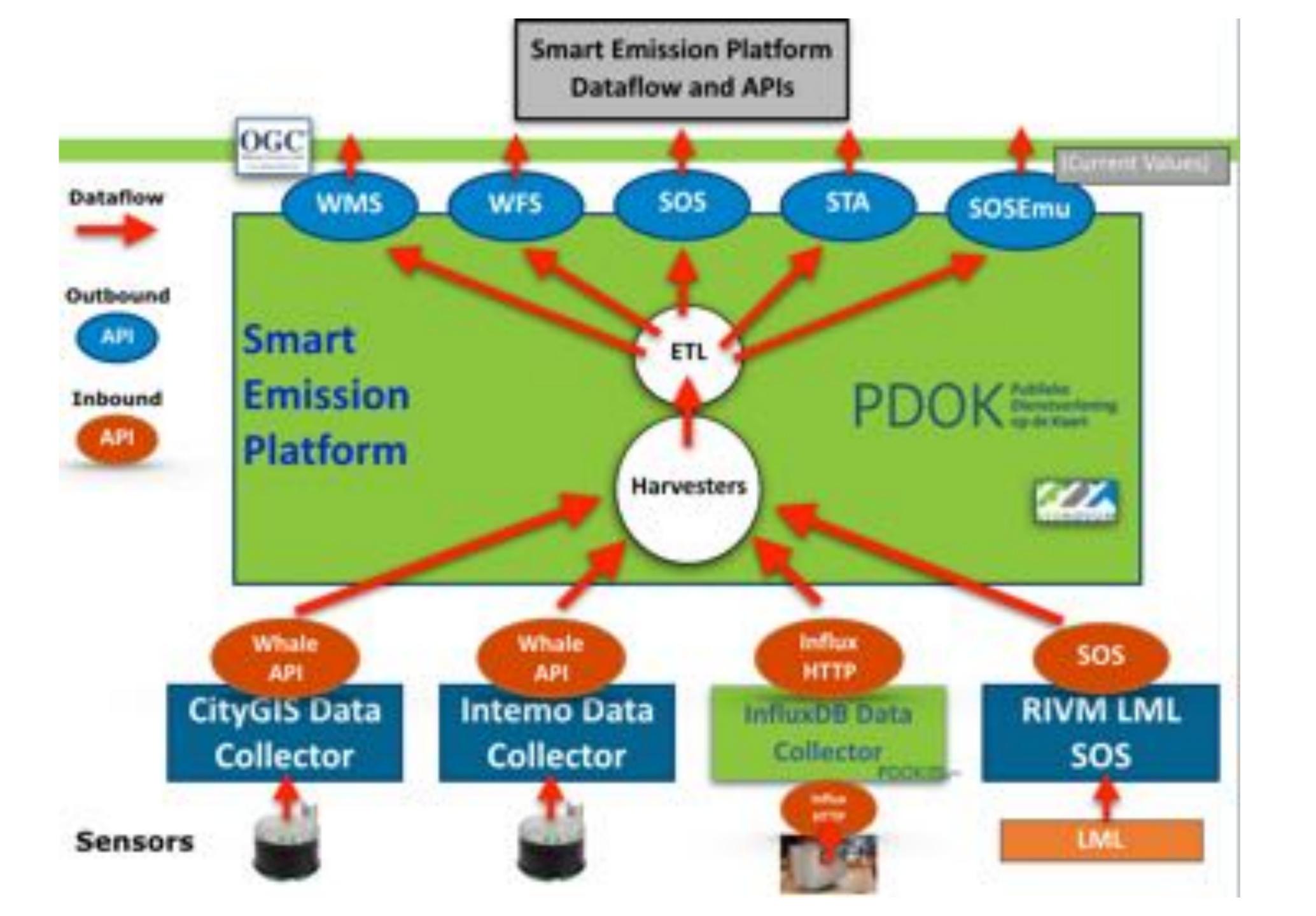
Smart Emission Platform Summary

- 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

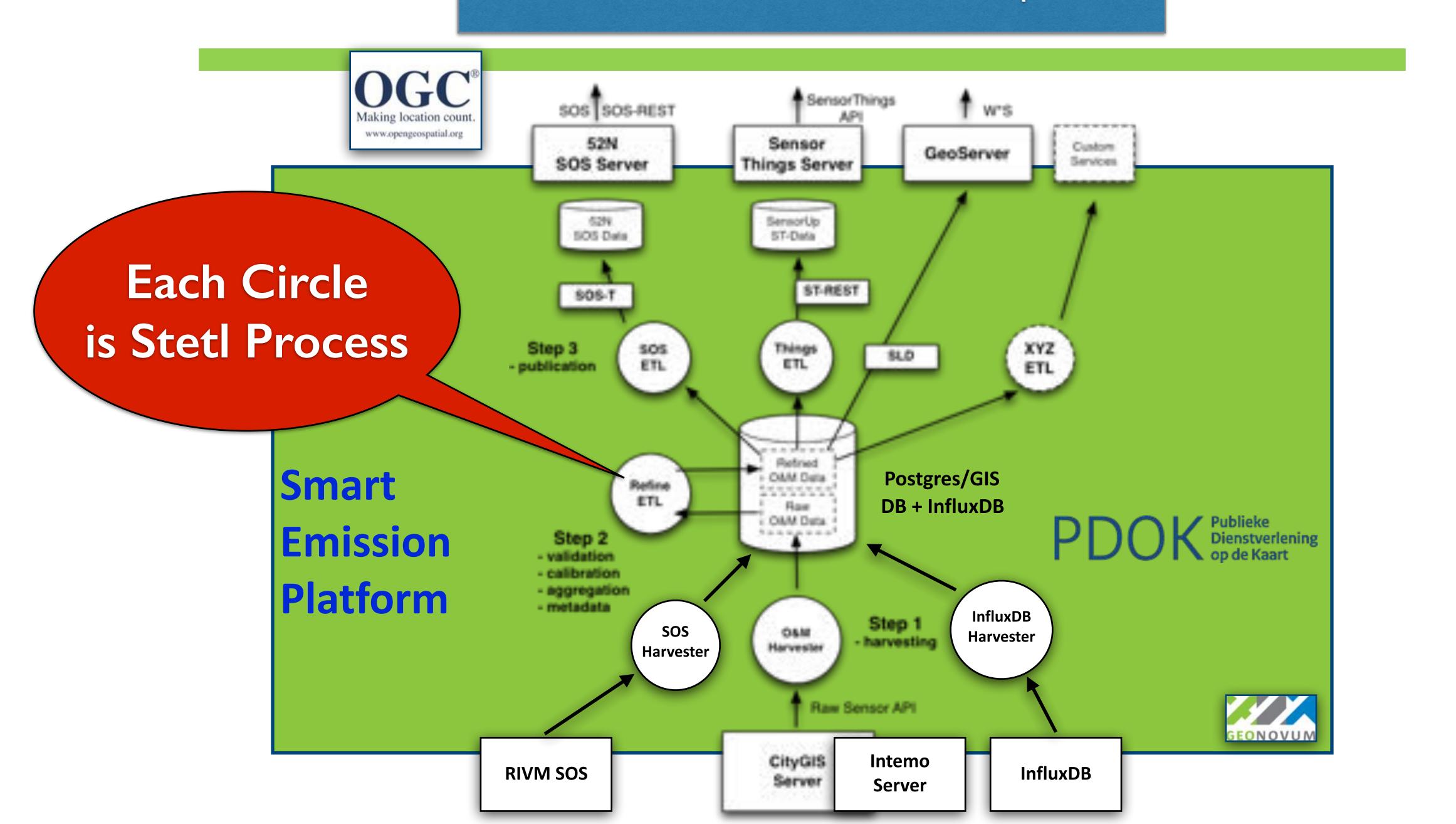
All ETL

(Harvesting, Validation, Calibration, Aggregation, Publication)

done with Stetl



Data Architecture with multi-Step ETL

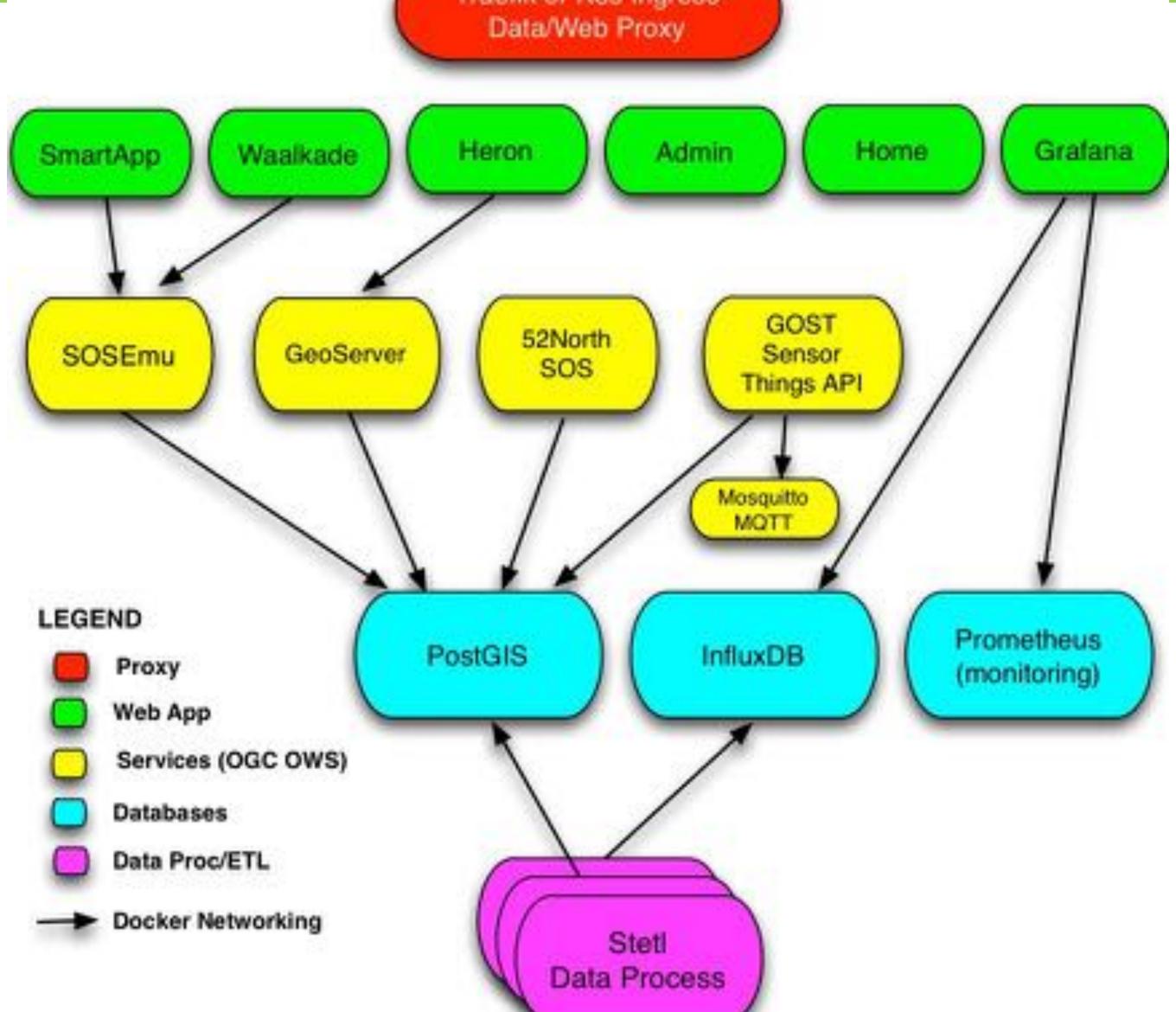




Smart Emission Docker Deployment



Traefik or K8s Ingress Data/Web Proxy



1. Dockerfile

```
FROM geopython/stetl:1.2
# One the standard Stetl Docker Tmage with some additional packages
LABEL maintainer- Just van den Broecke <justb#@gmail.com>
ENV ADD FYTHOW DEB PACKAGES- python-dov python-scipy python-seaborn python-matplotlib" \
    ADD PTTECH PIP PACKAGES- whool Goohash influxeb scikit-learn-0.18"
# Not now see! https://github.com/pyps/pip/issues/5240
# NUN pip install -- upgrade pip
NUM apt-get update &s apt-get --no-install-recommends install -y \
        $(ADD PETROW DEB FACKAGES) \
    as pip install $(ADD_PETNOW_PIP_PACKAGES) \
    44 apt-get purge -y python-dev \
    44 apt autoremove -y )
    as rm -rf /var/lib/apt/lists/*
# Copy relevant Files to work dir
ADD config /work/config
ADD smartem /work/smartem
ADD options/example.args /work/options/default.args
ADD scripts /work/scripts
NORKDIR /WOEK
ENTRYPOINT ["/work/scripts/entry.sh"]
```

3. Run ETL Process

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

2. entry.sh

```
#1/bin/bash
       # Generic ETL entry: override with specific ETL config and args file.
       WORK DIR-/work
       ed ${WORK_DIR}
       # Shorthand
       function log() (
           echo "entry.sh: $1"
       function orrer() {
         10g #
           exit -1
       # ETL Process can be specified via eav var or argument
       if [ "5(ETL_PROCESS)x" = "x" ];
           ETL PHOCESS" $1"
       # ETL Process can be specified via eav var or argument
       if ( "5(ETL PROCESS) = " a" 1;
           error "Error: so ETL process specified"
       BTL CONFIC FILE-config/8(ETL PROCESS).cfg
       if ( 1 -f $(ETL CONFIC FILE) ]
           error 'Error: cannot find ETL config file: 5(ETL CONFIG FILE)'
       else
           Jog "OK using config file: S(ETL COMFIG FILE) ... "
41
       # Start Stotl with config file
42
       export PTTHOMPATH-$ (WORK DIR) : $ (PYTHOMPATH)
4.3
       stetl -c $(ETL CONFIG FILE) -a options/default.args
```

Thank You!

www.stetl.org

github.com/geopython/stetl