

# The datahub

## De/blending museum data



Data has a better idea

## Setting the stage

In which I'll describe where we came from

## The Datahub Project

In which I'll show you an aggregation architecture

## The story thus far

In which I'll discuss the construction process

## What we learned

In which I'll conclude with a few take aways

**flemishartcollection**  
MUSEUMS OF FINE ARTS ANTWERP BRUGES GHENT LEUVEN OSTEND



KONINKLIJK  
MUSEUM  
VOOR SCHONE  
KUNSTEN  
ANTWERPEN

B R U  
G G E

MUSEA  
BRUGGE







Biographies

Collection

Research

Experience more

Home > Collection

Search in the complete collection



Filter by Museum

Royal Museum of Fine Arts Antwerp (155)

Groeninge Museum Bruges (105)

Museum of Fine Arts Ghent (73)

Saint Bavo Cathedral Ghent (29)

Museum M Leuven (23)

Museum Mayer van den Bergh Antwerp (20)

Sint-Janshospitaal Bruges (19)

Saint Salvator Cathedral Bruges (12)

OCMW Antwerpen (2)

Search

Enter your keywords:

Search

Search results



## Multiple organisations

Different local traditions, thesauri, various cataloguing rules (SPECTRUM), organisational contexts,...

## Multiple registration systems

TMS, Adlib, CollectiveAccess, closed/open source, Lack of API's, non standardised API's,...

## Multiple end user applications

various websites, historically grown, different contractors, various CMS systems, different ways to deliver data,...

# Manual exchange

## Different ways

Excel, CSV, vendor formats. WeTransfer, e-mail,...

## Error prone

Corrupt exports, wrong data exported, wrong version passed on, stuff gets lost along the way,...

## High overhead costs

Time and money (communication, \$/hour)

## High latency

What's online is not really up to date

# Herding cats





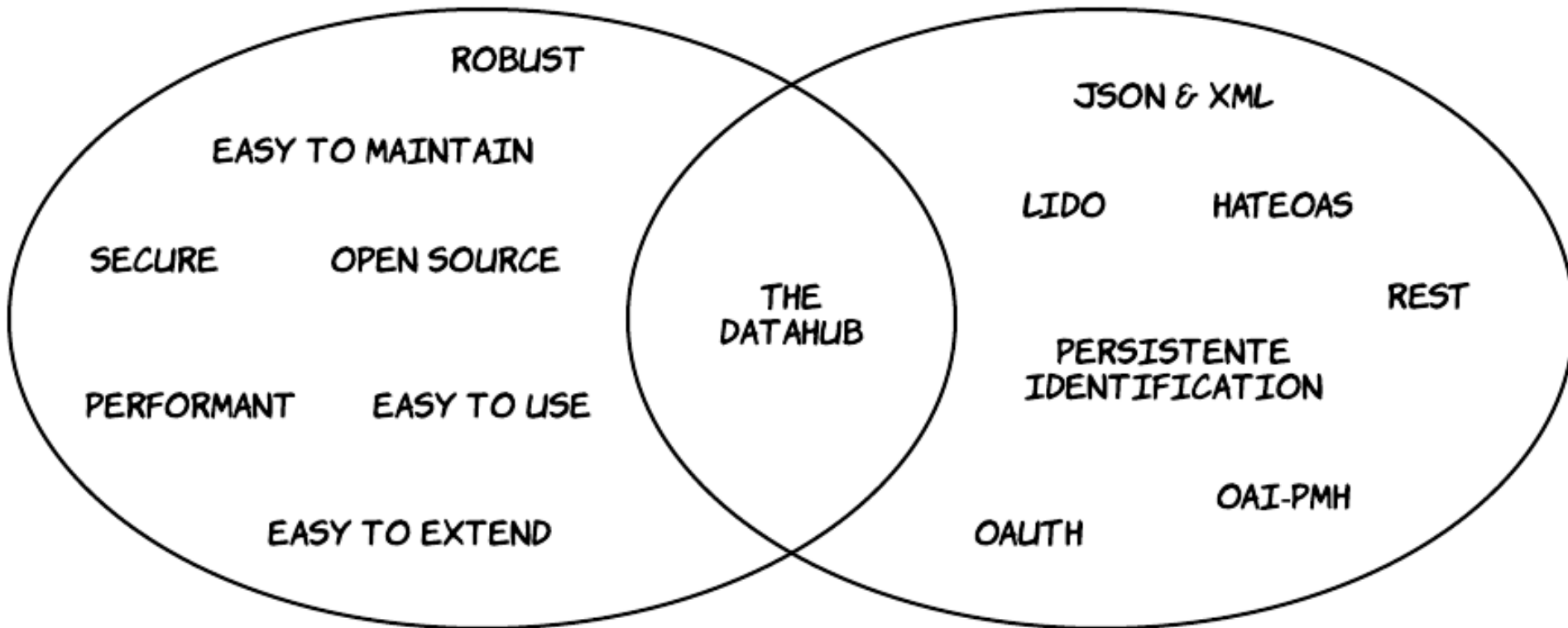
# A modern ecosystem



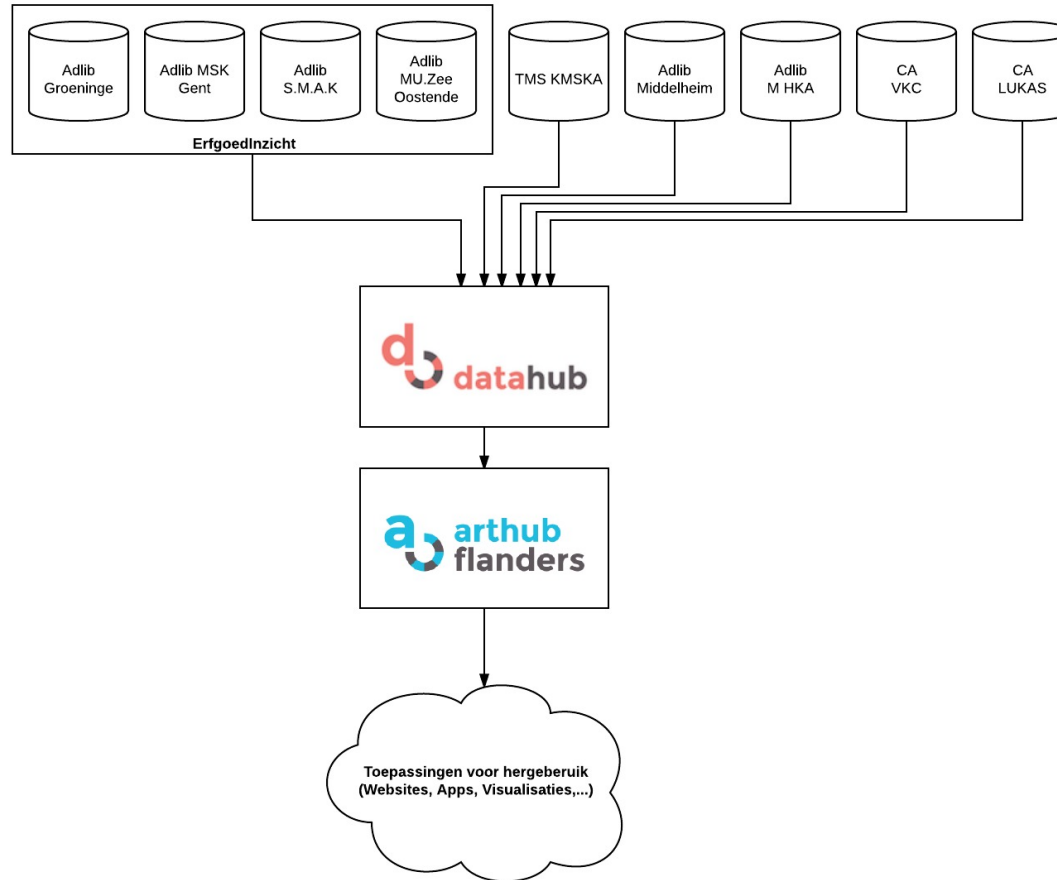
# The Datahub Project



# Aggregator



# Local aggregator



# Arthub Flanders

This datahub currently contains 15629 records. These records are published through a [REST API endpoint](#) and an [OAI-PMH endpoint](#).

This datahub is managed by [Vlaamse Kunstcollectie vzw](#). Reach out via e-mail at [noreply@datahub.inuits.eu](mailto:noreply@datahub.inuits.eu).

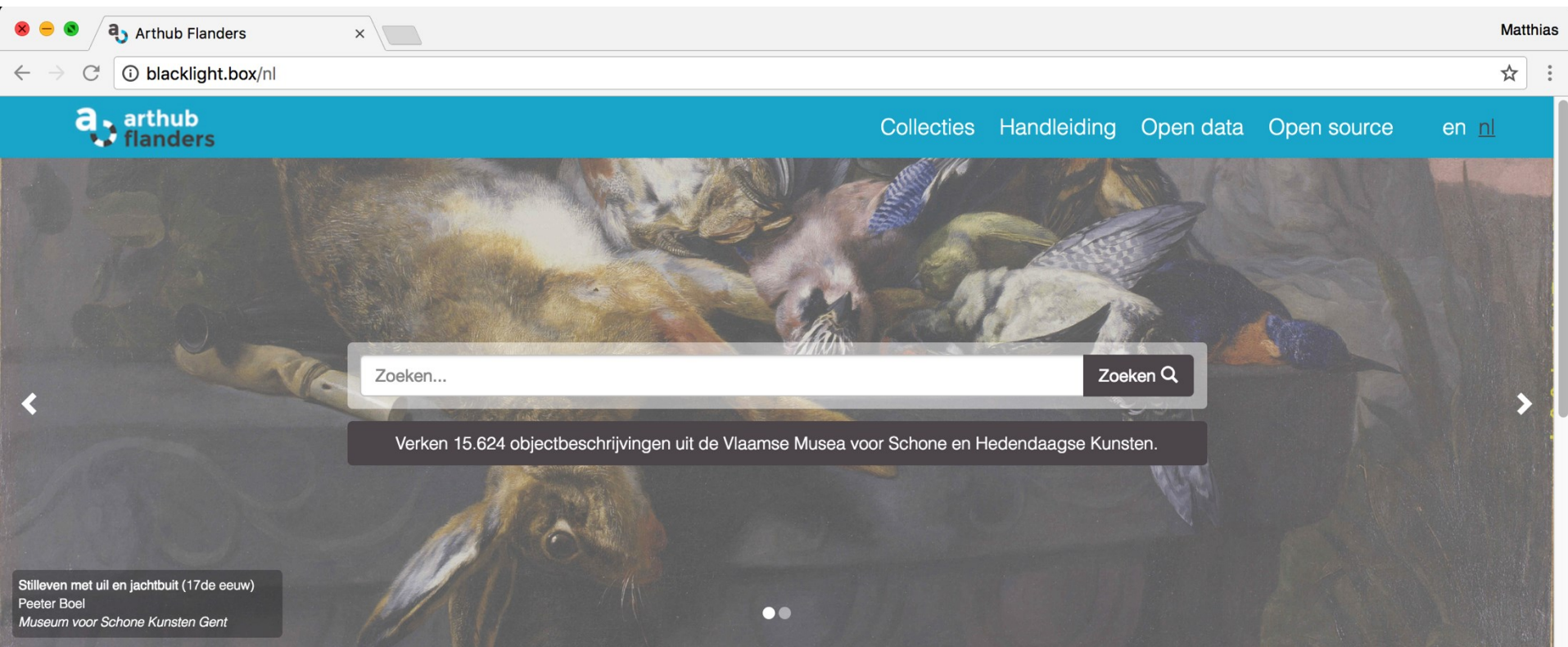


datahub.box/oai/?metadataPrefix=oai\_lido&verb=ListRecords

This XML file does not appear to have any style information associated with it. The document tree is shown below.

```
<?xml version="1.0" encoding="UTF-8" standalone="yes"?>
<OAI-PMH xmlns="http://www.openarchives.org/OAI/2.0/" xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
  xsi:schemaLocation="http://www.openarchives.org/OAI/2.0/ http://www.openarchives.org/OAI/2.0/OAI-PMH.xsd">
  <responseDate>2018-06-01T13:15:31Z</responseDate>
  <request metadataPrefix="oai_lido" verb="ListRecords">http://datahub.box/</request>
  <ListRecords>
    <record>
      <header>
        <identifier>
          oai:datahub.vlaamsekunstcollectie.be:groeningemuseum.be:0000_GRO1561_I
        </identifier>
        <timestamp>2018-05-02T14:42:04Z</timestamp>
      </header>
      <metadata>
        <lido:lido xmlns:gml="http://www.opengis.net/gml" xmlns:lido="http://www.lido-schema.org" xmlns:xlink="http://www.w3.org/1999/xlink">
          <lido:lidoRecID lido:pref="alternate" lido:type="purl" lido:source="Musea Brugge - Groeningemuseum" lido:label="dataPID">
            http://groeningemuseum.be/collection/work/data/0000_GRO1561_I
          </lido:lidoRecID>
          <lido:lidoRecID lido:pref="preferred" lido:type="urn" lido:source="Vlaamse Kunstcollectie - Arthub Flanders" lido:label="dataPID">
            oai:datahub.vlaamsekunstcollectie.be:groeningemuseum.be:0000_GRO1561_I
          </lido:lidoRecID>
          <lido:objectPublishedID lido:type="purl" lido:source="Musea Brugge - Groeningemuseum" lido:label="workPID">
            http://groeningemuseum.be/collection/work/id/0000_GRO1561_I
          </lido:objectPublishedID>
          <lido:category>
            <lido:conceptID lido:type="purl" lido:source="cidoc-crm">http://www.cidoc-crm.org/crm-concepts/E22</lido:conceptID>
            <lido:term>Man-Made Object</lido:term>
          </lido:category>
          <lido:descriptiveMetadata xml:lang="nl">
            <lido:objectClassificationWrap>
              <lido:objectWorkTypeWrap>
                <lido:objectWorkType>
                  <lido:conceptID lido:pref="preferred" lido:type="local" lido:source="Adlib">20000001</lido:conceptID>
                  <lido:conceptID lido:pref="alternate" lido:type="purl" lido:source="AAT">http://vocab.getty.edu/aat/300033618</lido:conceptID>
                  <lido:term lido:pref="preferred" xml:lang="nl">schilderingen</lido:term>
                  <lido:term lido:pref="alternate" xml:lang="nl">schilderingen</lido:term>
                </lido:objectWorkType>
              </lido:objectWorkTypeWrap>
            </lido:descriptiveMetadata>
            <lido:classification>
              <lido:conceptID lido:pref="preferred" lido:type="local" lido:source="Adlib">20000153</lido:conceptID>
```

```
1 // 20180601151621
2 // http://datahub.box/api/v1/data.json
3
4 {
5   "offset": 0,
6   "limit": 5,
7   "total": 15629,
8   "_links": {
9     "self": {
10      "href": "/api/v1/data?limit=5"
11    },
12    "first": {
13      "href": "/api/v1/data?limit=5"
14    },
15    "last": {
16      "href": "/api/v1/data?offset=15625&limit=5"
17    },
18    "next": {
19      "href": "/api/v1/data?offset=5&limit=5"
20    }
21  },
22  "_embedded": {
23    "records": Array[5][
24      {
25        "id": "5ae9ce3c72a84303de6f2ada",
26        "created": "2018-05-02T09:42:04-05:00",
27        "updated": "2018-05-02T09:42:04-05:00",
28        "json": Array[6][
29          {
30            "name": "{http://www.lido-schema.org}lidoRecID",
31            "value": "http://groeningemuseum.be/collection/work/data/0000_GR01561_I",
32            "attributes": {
33              "{http://www.lido-schema.org}pref": "alternate",
34              "{http://www.lido-schema.org}type": "purl",
35              "{http://www.lido-schema.org}source": "Museum Brugge - Groeningemuseum",
36              "{http://www.lido-schema.org}label": "dataPID"
```



## Wat is Arthub Flanders?

Arthub Flanders verzamelt beschrijvingen over kunst- en erfgoedobjecten opgesteld en beheerd door de Vlaamse musea voor Schone en Hedendaagse Kunsten. Arthub Flanders publiceert deze beschrijvingen in open formaten en onder een open licentie zodat iedereen ze kan hergebruiken in eigen toepassingen

Alle velden

Zoeken...

Zoeken

## Verfijn uw zoekopdracht

Periode



Instelling



Type



Subtype



Materiaal



Onderwerp



« Vorige | 1 - 10 van 15.624 | Volgende »

Sorteer op relevantie

10 per pagina

## 1. Johannes predikt tot de menigte



Vervaardiger: [kopie naar Bruegel, Pieter I](#)  
[atelier van Brueghel, Pieter II](#)  
[atelier van Brueghel, Jan I](#)

Periode: [17de eeuw](#)

Instelling: [Musea Brugge - Groeningemuseum](#)

Type: [schilderingen](#)

Onderwerp: [religieuze voorstellingen](#) en [landschappen](#)

Data PID: [http://groeningemuseum.be/collection/work/data/0000\\_GRO1561\\_I](http://groeningemuseum.be/collection/work/data/0000_GRO1561_I)

## 2. Opvoeding van Maria



Vervaardiger: [Garemijn, Jan Anton](#)

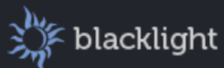
Periode: [18de eeuw](#)

Instelling: [Musea Brugge - Groeningemuseum](#)

Type: [schilderingen](#)

Onderwerp: [religieuze voorstellingen](#)

Data PID: [http://groeningemuseum.be/collection/work/data/0000\\_GRO1561\\_I](http://groeningemuseum.be/collection/work/data/0000_GRO1561_I)

[GITHUB](#)[WIKI](#)[EVENTS](#)[RELEASE NOTES](#)[APPS & DEMOS](#)

# Blacklight

A multi-institutional open-source collaboration building a better discovery platform framework

[Learn how to get started](#)[Examples](#)

## Featured Plugins

### Blacklight MARC

Library catalog enhancements for Blacklight.

### Spotlight

Enable librarians, curators, and others who are responsible for digital collections to create attractive, feature-rich websites that highlight these collections.

### GeoBlacklight

A multi-institutional open-source collaboration building a better way to find and share geospatial data.



# Catmandu

- the data processing toolkit -



## What is Catmandu?

Catmandu is a command line tool to access and convert data from your digital library, research services or any other open data sets.

## Features

```
[General]
id_path = '_metadata.object_number'

[Importer]
plugin = EIZ

[plugin_importer_EIZ]
endpoint = https://endpoint.org
metadata_prefix = oai_adlib
username = username
password = secretpassword
handler = adlib

[Fixer]
plugin = Fix

[plugin_fixer_Fix]
condition_path = "_metadata.institution\\.name.value"
fixers = MSK, GRO

[plugin_fixer_MSK]
condition = 'Museum voor Schone Kunsten Gent'
file_name = '/Users/matthiasvandermaesen/Workspace/Datahub-Fixes/msk_oai_adlib.fix'

[plugin_fixer_GRO]
condition = 'Musea Brugge - Groeningemuseum'
file_name = '/Users/matthiasvandermaesen/Workspace/Datahub-Fixes/groeninge_oai_adlib.fix'

[Exporter]
plugin = Arthub

[plugin_exporter_Arthub]
datahub_url = http://datahub.box
datahub_format = LIDO
oauth_client_id = slightlylesssecretpublicid
oauth_client_secret = supersecretsecretphrase
oauth_username = admin
oauth_password = datahub
```



matthiasvandermaesen at Artemis in ~

\$ dhconveyor

commands: list the application's commands

help: display a command's help screen

index: Transport data from a flat file to a data index in bulk.

transport: Transport data from a data source to a data sink.

matthiasvandermaesen at Artemis in ~

\$ dhconveyor transport -p ~/Workspace/arthur-pipelines/erfgoedinzicht.ini -v

Loading pipeline configuration...

Initializing importer/exporter...

Initializing fixers...

Importing data from source...

✓ - Item #1 : 0000.GR01561.I (id): exported.

✓ - Item #2 : 0000.GR01390.I (id): exported.

✓ - Item #3 : 0000.GR00128.I (id): exported.

✓ - Item #4 : 0000.GR01476.I (id): exported.

✓ - Item #5 : 0000.GR00479.I (id): exported.

✓ - Item #6 : 0000.GR01372.I (id): exported.

✓ - Item #7 : 0000.GR01360.I (id): exported.

✓ - Item #8 : 0000.GR01359.I (id): exported.

✓ - Item #9 : 0000.GR00227.I (id): exported.

✓ - Item #10 : 0000.GR01280.I (id): exported.

✓ - Item #11 : 0000.GR01243.I (id): exported.

✓ - Item #12 : 0000.GR00299.I (id): exported.

✓ - Item #13 : 0000.GR01230.I (id): exported.

^C

matthiasvandermaesen at Artemis in ~

\$

Resolver

Entities

Users

Stats

Settings

Import & Export

Sign out

PID: 0000\_GRO0018\_I

Add entity ...

Edit entity ...

Documents ...

Persistent URIs ...

[http://resolver.vlaamsekunstcollectie.be/collection/0000\\_GRO0018\\_I](http://resolver.vlaamsekunstcollectie.be/collection/0000_GRO0018_I)

[http://resolver.vlaamsekunstcollectie.be/collection/0000\\_GRO0018\\_I/saint-luke-painting-the-madonna](http://resolver.vlaamsekunstcollectie.be/collection/0000_GRO0018_I/saint-luke-painting-the-madonna)

[http://resolver.vlaamsekunstcollectie.be/collection/work/data/0000\\_GRO0018\\_I/html](http://resolver.vlaamsekunstcollectie.be/collection/work/data/0000_GRO0018_I/html)

[http://resolver.vlaamsekunstcollectie.be/collection/work/data/0000\\_GRO0018\\_I/html/saint-luke-painting-the-madonna](http://resolver.vlaamsekunstcollectie.be/collection/work/data/0000_GRO0018_I/html/saint-luke-painting-the-madonna)

[http://resolver.vlaamsekunstcollectie.be/collection/work/data/0000\\_GRO0018\\_I](http://resolver.vlaamsekunstcollectie.be/collection/work/data/0000_GRO0018_I)

[http://resolver.vlaamsekunstcollectie.be/collection/work/representation/0000\\_GRO0018\\_I/1](http://resolver.vlaamsekunstcollectie.be/collection/work/representation/0000_GRO0018_I/1)

[http://resolver.vlaamsekunstcollectie.be/collection/work/representation/0000\\_GRO0018\\_I/1/saint-luke-painting-the-madonna](http://resolver.vlaamsekunstcollectie.be/collection/work/representation/0000_GRO0018_I/1/saint-luke-painting-the-madonna)



# The story thus far



# Assumptions / Reality

- We had a fixed, limited budget
  - Estimated timeline 3 to 6 months.
  - A production ready version.
- 
- Contractor delivered a prototype version.
  - We over-extended the timing.
  - Switch to DIY development after 6 months.
  - Scope changes as we went along.

# What happened?

- We underestimated the ETL workload
- We overestimated contractor engagement
- We underestimated organisational complexity

# Wicked ETL

- Context really matters
- Getting intimate with the domain takes time
- Integrating data across network is challenging.
- Difficult to guestimate complexity up front

# Wicked ETL

Context really matters

- Machines  
Legacy software, lack of infrastructure,...
- People  
Data means nothing until it gets interpreted.  
But, different perceptions of reality...
- Content  
Driven by tradition, software, people.



# Wicked ETL

Data modelling is a wicked challenge

- Mapping to standardised exchange formats  
... and their specific data models
- Normalisation and enrichment  
... are we taking about the same thing?
- Context specific concerns  
... Copyright, privacy, security, authority

# Procurement

- Build-to-print vs build-to-spec.
- You outsource the process, not the project.
- Is contractor service a good fit?
- Relationship with the contractor!
- Procurement is part of the design process

# DIY development

- Knowledge domain and technical experience
- Flexibility to build exactly what you need
- Reduces dependency on a specific contractor
- Requires in-house competences
- Payroll is a hidden cost
- The Bus Factor risk

# Lessons learned



## Own your project

Define the project process you're going to follow

## Actively involve your stakeholders

Challenge your own assumptions, but keep your focus!

## Actively be involved in the process

Don't assume a vendor will solve things for you.

# Be mindful about the budget

Fixed price vs Fixed budget

# In source talented specialists you need

Identify right profile: IA, Dev, PM, UX,...

# Outsource placing the kitchen sink

Stock off-the-shelf website or web app

## Document all the things

Be mindful about the human who comes after you!

## Don't do elaborate specifications up front

Nobody is interested in paper tigers.

## Make your hands dirty

Try tools up front. Identify the big hurdles early.



# Thank you!

<https://github.com/thedatahub>

<https://thedatahub.github.io>

<http://www.flemishartcollection.be>

T: @netsensei