



národní
úložiště
šedé
literatury

The datahub

Vandermaesen, Matthias
2018

Dostupný z <http://www.nusl.cz/ntk/nusl-403475>

Dílo je chráněno podle autorského zákona č. 121/2000 Sb.

Licence Creative Commons Uveďte původ-Neužívejte komerčně-Nezpracovávejte 4.0

Tento dokument byl stažen z Národního úložiště šedé literatury (NUŠL).

Datum stažení: 19.04.2024

Další dokumenty můžete najít prostřednictvím vyhledávacího rozhraní nusl.cz .

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)

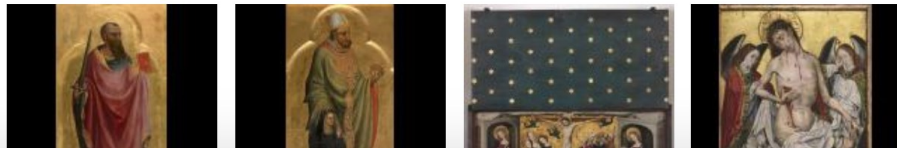
OCMW Antwerpen (2)

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



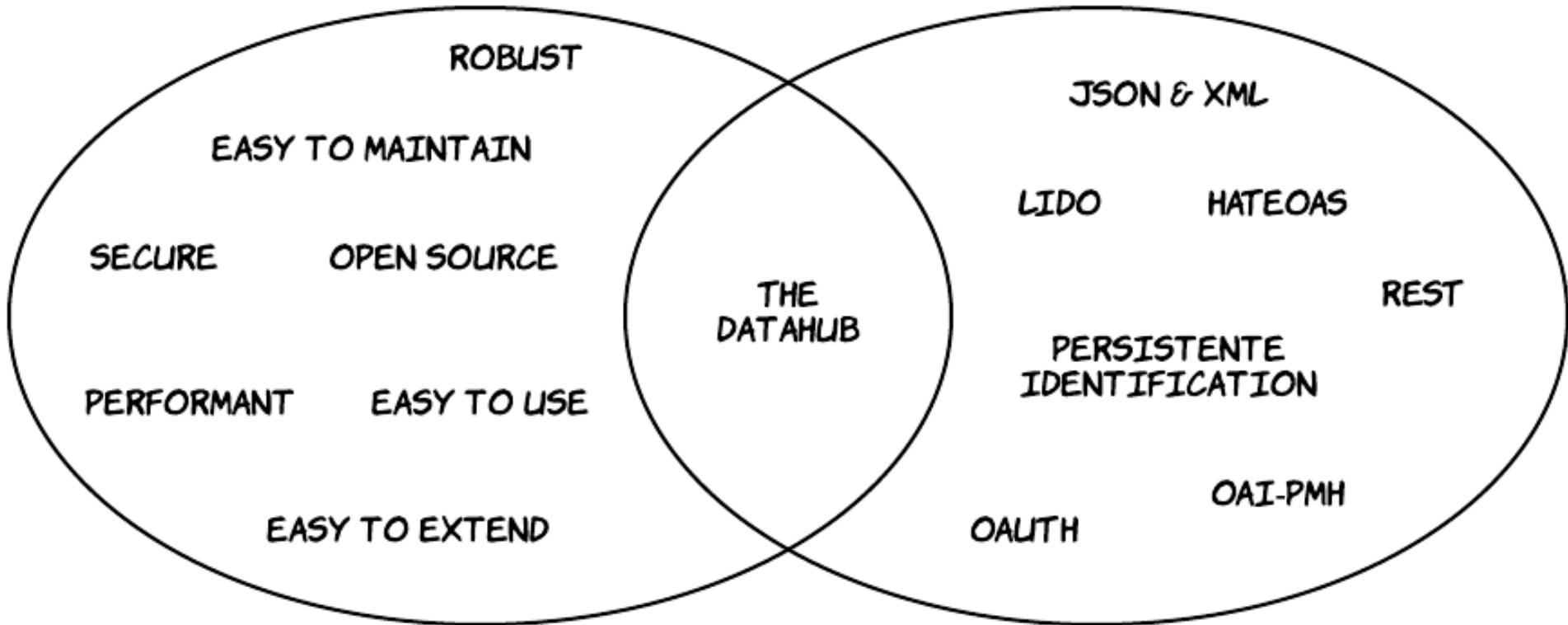
A modern ecosystem



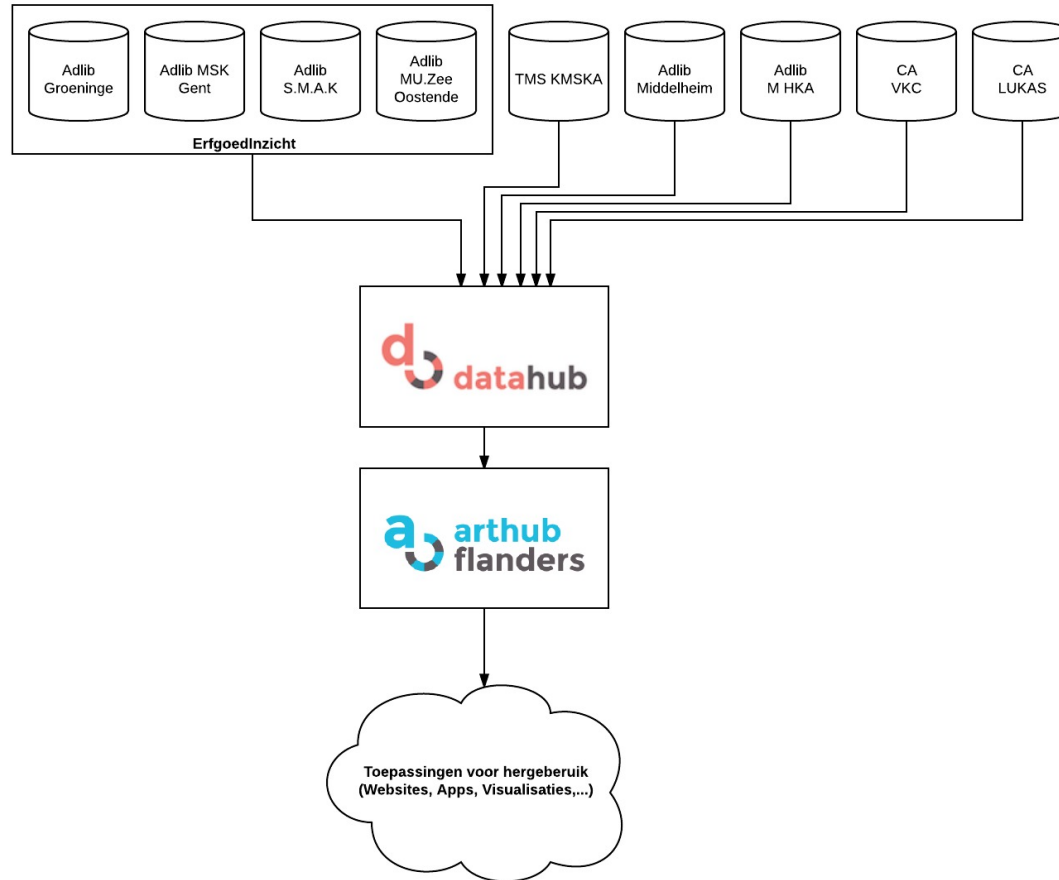
The Datahub Project



Aggregator



Local aggregator



Arthub Flanders

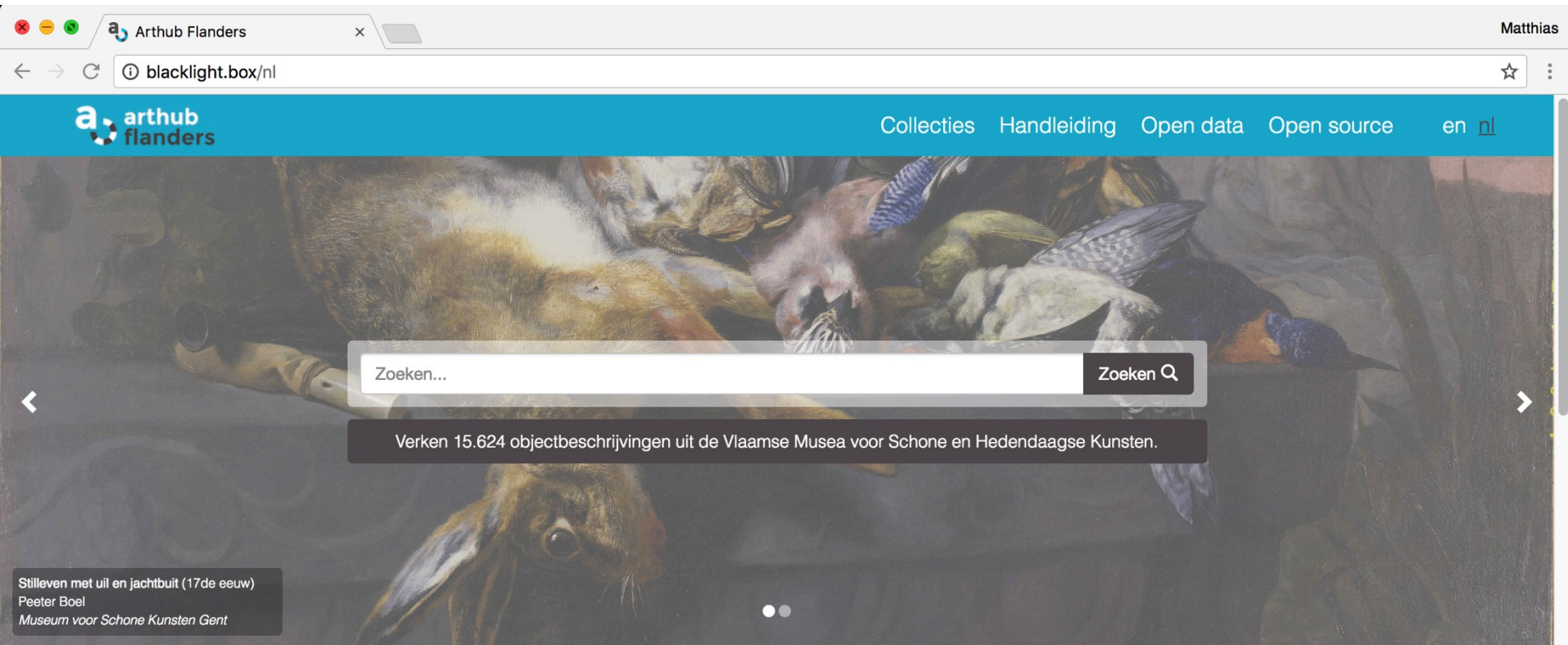
This datahub currently contains 15629 records. These records are published through a [REST API endpoint](#) and an [OAI-PMH endpoint](#).

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

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" ?>
<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:objectClassificationWrap>
            <lido:classificationWrap>
              <lido:classification>
                <lido:conceptID lido:pref="preferred" lido:type="local" lido:source="Adlib">20000153</lido:conceptID>
              </lido:classification>
            </lido:classificationWrap>
          </lido:descriptiveMetadata>
        </lido:lido>
      </metadata>
    </record>
  </ListRecords>
</OAI-PMH>
```

```
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_GRO1561_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": "Musea 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

2. [Opvoeding van Maria](#)

Vervaardiger: [Garemijn, Jan Anton](#)

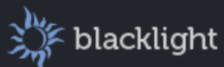
Periode: [18de eeuw](#)

Instelling: [Musea Brugge - Groeningemuseum](#)

Type: [schilderingen](#)

Onderwerp: [religieuze voorstellingen](#)

Data PID: http://groeningemuseum.be/collection/work/data/0000_GRO1561_I



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

groeninge_oai_adlib.fix x

```
83 ### LIDO lidoRecID
84
85 # ID
86 #
87 # The ID in Solr is based on the data_pid. The data_pid is converted to a string
88 # which can be safely used as an identifier in Project Blacklight. The format of
89 # the ID field looks like this:
90 #
91 # oai:datahub.vlaamsekunstcollectie.be:<domain>:<identifier>
92 # ex. oai:datahub.vlaamsekunstcollectie.be:kmsa.be:254
93 # ex. oai:datahub.vlaamsekunstcollectie.be:collectievlaamsegemeenschap.be:837
94 #
95 # Note: the .tld is stripped from the domainname because the . (dot) breaks the
96 # route matching algorithm.
97
98
99 unless is_array(or_record.dataPid)
100   move_field(or_record.dataPid, or_record.tmp)
101   set_array(or_record.dataPid)
102   move_field(or_record.tmp, or_record.dataPid.$last)
103 end
104
105 do list(path:or_record.dataPid, var: c)
106
107   if all_match('c.digital_reference\description.value', 'datapid')
108
109     lido_baseid(
110       lidoRecID,
111       c.digital_reference,
112       -type: purl,
113       -source: 'Musea Brugge - Groeningemuseum',
114       -label: dataPID,
115       -pref: alternate
116     )
117
118     copy_field(c.digital_reference, or_record.oaiPid)
119     parse_text(or_record.oaiPid, '.*://([A-Za-z0-9-\.]+)/collection/work/data/(.*)')
120     join_field('or_record.oaiPid', ':')
121     prepend('or_record.oaiPid', 'oai:datahub.vlaamsekunstcollectie.be:')
122
123     lido_baseid(
124       lidoRecID,
125       or_record.oaiPid,
```

```

126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
```

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/arthub-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_Ihttp://resolver.vlaamsekunstcollectie.be/collection/0000_GRO0018_I/saint-luke-painting-the-madonnahttp://resolver.vlaamsekunstcollectie.be/collection/work/data/0000_GRO0018_I/htmlhttp://resolver.vlaamsekunstcollectie.be/collection/work/data/0000_GRO0018_I/html/saint-luke-painting-the-madonnahttp://resolver.vlaamsekunstcollectie.be/collection/work/data/0000_GRO0018_Ihttp://resolver.vlaamsekunstcollectie.be/collection/work/representation/0000_GRO0018_I/1http://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 guesstimate 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