

Improving data quality

A Danish case study

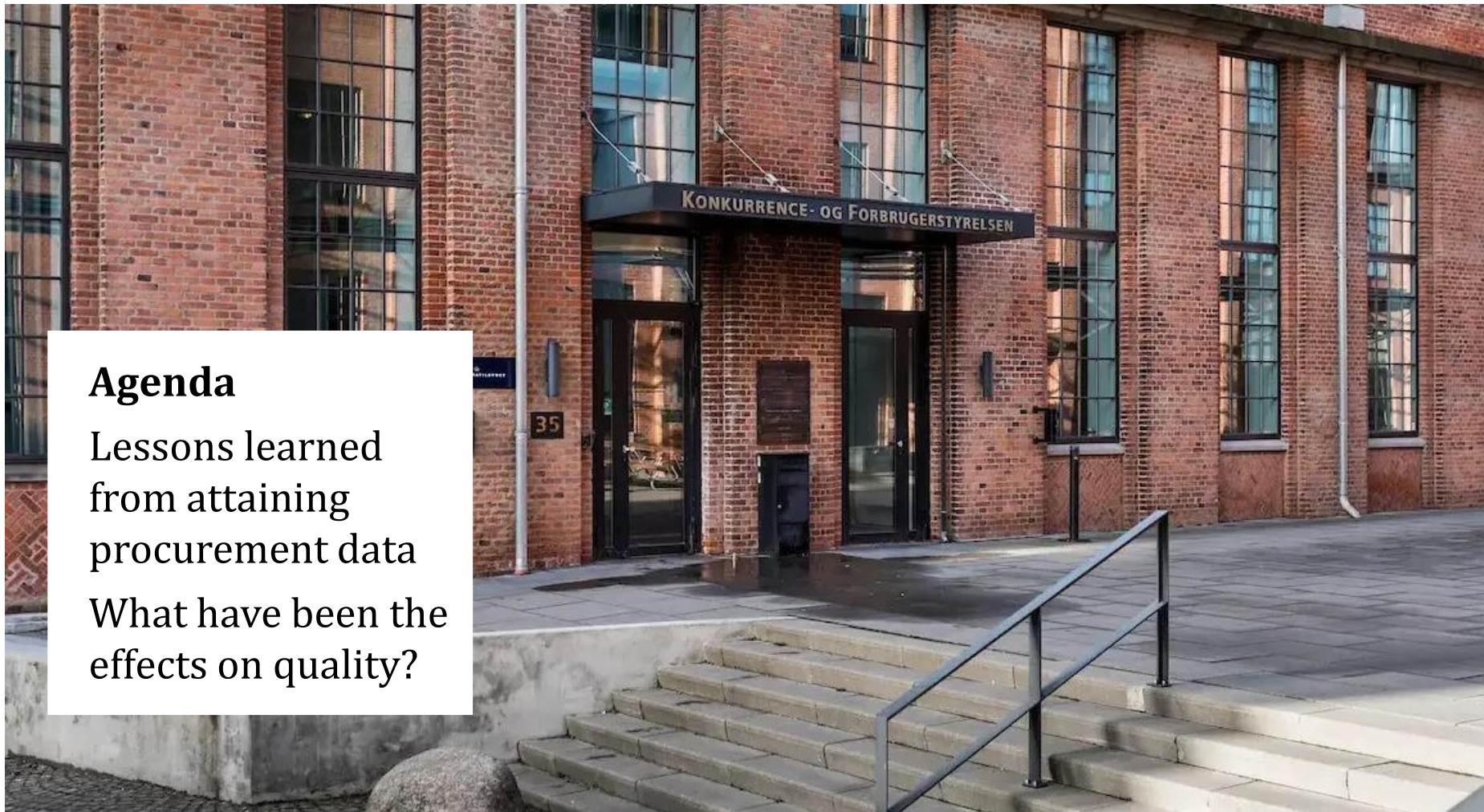
Lily Jacobsen Persson – 3. December 2024



Agenda

Lessons learned
from attaining
procurement data

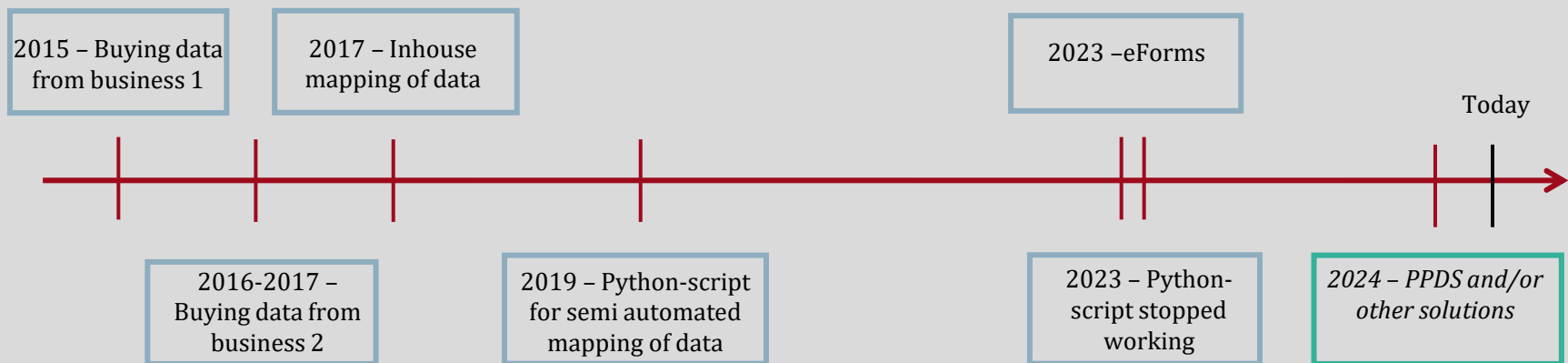
What have been the
effects on quality?



Lessons learned from attaining procurement data



Timeline of procurement data



2015 - Buying data from business 1

Purpose and approach

- Measure the impact of the new directive
- Very manual approach centered around CANs and associated documents (the CN, modifications, prior notices)

Lessons learned

- A lot of back and forth with the supplier on getting the desired dataset
- It cost a bit



2017 – Buying data from business 2

Purpose and approach

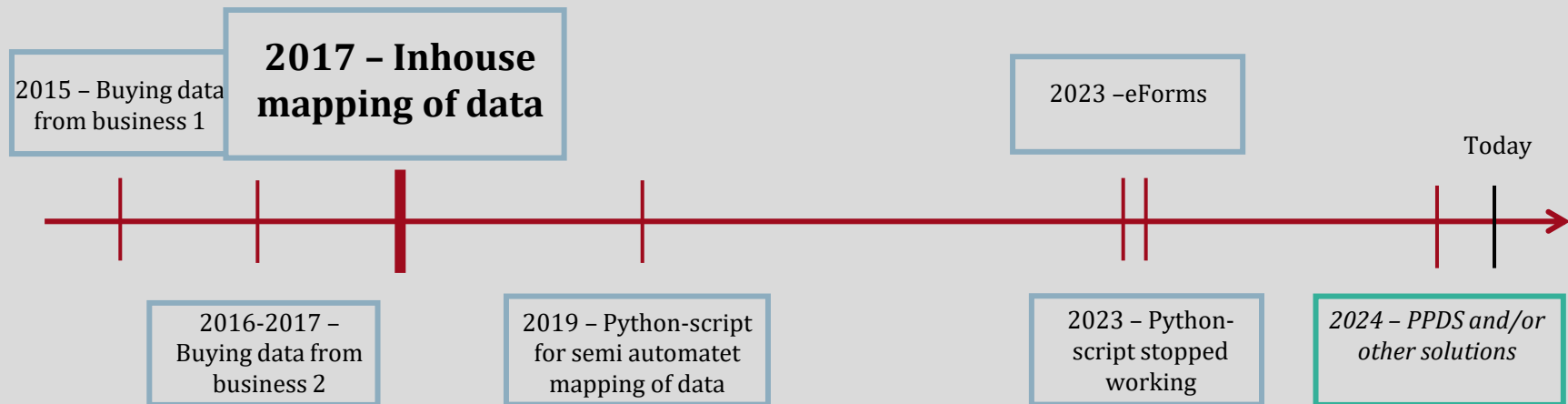
- A new company won the contract mapping 2016 and half of 2017
- Different approach with web-scraping TED to obtain data

Lessons learned

- Lower data quality and harder to identify the errors due to scripting errors
- No ongoing information on procurement
- Lack of control over the process and decisions
- Still cost a bit



Timeline of procurement data



2017 – Our own manual mapping

Purpose and approach

- New in-house unit with 3-4 student assistants
- Ongoing discussions on good practices of mapping data, methodology, practises/reality vs. rules with our legal experts
- Procedures for contacting procuring authorities when identifying missing information or errors
- Procedures for quality assurance every month



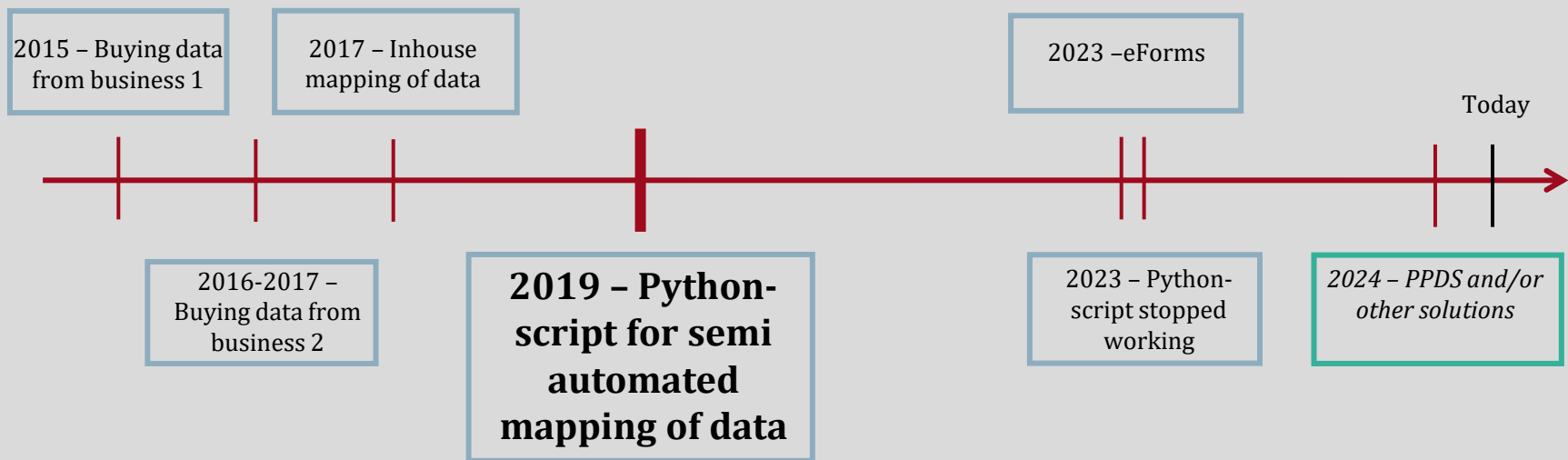
2017 – Our own manual mapping

Lessons learned

- Knowledge on what information can be mapped and how to map tenders
- Understanding practises of procuring authorities
- Limitations of data
- Less expensive and other task for student assistants
- Not efficient – and hard to keep the task interesting

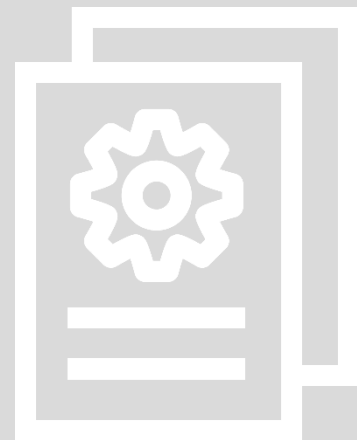


Timeline of procurement data



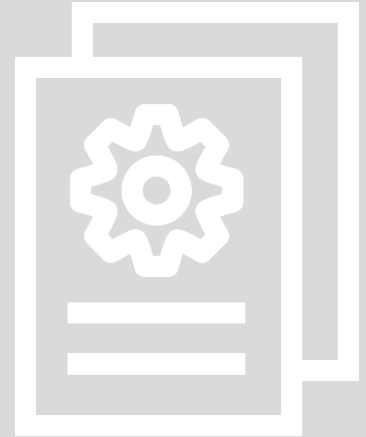
2019 – Automation of mapping data

- We had a well-developed methodology that could be translated to code
- A student assistant with a gift for programming set forth creating a web-scraping script



2019 – Automation of mapping data

- Built a web-scraping script based on our methodology
- Mapped data from CANs and associated notices
- Corrected information automatically (when possible)
- Highlighted/marked fields and possible errors to be checked by a student assistant



temp.py x Input_H_script.py x Bok_KFST.py x opdelKontrakter.py x

```

660     except AttributeError:
661         udbudDirektiv = ''
662         pass
663
664     """
665     Først og fremmest undersøges, hvor mange delkontrakter, der er.
666     Da den variabel, som Ordregiver kan bruge er dårlig, anvendes en anden tilgang,
667     hvor antallet af "Betegnelse" der inddeler en ny delkontrakt anvendes som indikator
668     Der anvendes samme tilgang som i kontrakttildelingssektionen
669     Vi ønsker, at finde span med id='id-II.' dens parent's parent er netop
670     Det II: Genstand
671     """
672     """ VARIABEL Antal delkontrakter """
673
674     """ VARIABEL Opdelt udbud """
675     delIIGenstandSuppe = boikReqSuppe.find("span", {"id": "id-II."}).parent.parent
676     delIIGenstand = delIIGenstandSuppe.getText()
677
678     # der fratrækkes 1 da den første Betegnelse altid er der uanset antal delkontrakter
679     Antal_delkontrakter_kortlagt = delIIGenstand.count("Betegnelse:") - 1
680     if Antal_delkontrakter_kortlagt == 1:
681         opdelt_udbud_BOIK = 'nej'
682
683     # hvis der er flere delkontrakter end 1, da er det opdelt
684     elif Antal_delkontrakter_kortlagt > 1:
685         opdelt_udbud_BOIK = 'ja'
686
687         """
688         Det sker sommetider, at "Betegnelse" ikke fremgår af afsnittet. Så rettes den dog til
689         """
690
691     elif Antal_delkontrakter_kortlagt < 1:
692         Kommentar = tilføjKommentar(Kommentar, 'Usikkerhed om antallet af delkontrakter på udbuddet. ', 'Vigtigt')
693         Antal_delkontrakter_kortlagt = str(1)
694         Antal_delkontrakter_kortlagt = tilføjKommentar(Antal_delkontrakter_kortlagt, ' ', 'Vigtigt')
695         opdelt_udbud_BOIK = 'nej'
696
697     # Antal delkontrakter kortlagt konverteres til string for at arbejde med den senere
698     Antal_delkontrakter_kortlagt = str(Antal_delkontrakter_kortlagt)
699
700     """
701     Her tjekkes om enten udbudsbeholdningsregulering eller BOIK siger opdelt udbud
702     og i så fald antages det, at der er tale om opdelt udbud.
703     """
704     if opdelt_udbud_BOIK == 'ja' or opdelt_udbud_Udbud == 'ja':
705         Opdelt_udbud = 'Ja'
706     else:
707         Opdelt_udbud = 'Nej'
708
709     # Laver dummy, der gør det muligt, at tilføje delkontraktadskiller. 1 når opdelt 0 ellers
710     if Opdelt_udbud == 'Ja':
711         DummyOpdelt = 1
712     else:
713         DummyOpdelt = 0
714
715     """
716     Rettelse af Antal_delkontrakter_kortlagt i tilfælde af at ordregiveren ikke har
717     udfyldt det igen på BOIK'en sørges for at, hvis antallet af kontrakttildelinger matcher antallet af
718     delkontrakter på udbuddet, da antages det at, der er ligeså mange delkontrakter på BOIK'en.
719     Der tilføjes dog også en kommentar.
720     """
721     spanKontrakttildelinger = boikReqSuppe.findAll("span", {"id": re.compile('id\d+-V.')}))
722

```

Source Console Object

Usage

Here you can get help of any object by pressing **Ctrl+I** in front of it, either on the Editor or the Console.

Help can also be shown automatically after writing a left parenthesis next to an object. You can activate this behavior in [Preferences > Help](#).

[New to Spyder? Read our tutorial](#)

[Help](#) [Variable Explorer](#) [Plots](#) [Files](#)

Console 1/A x

Python 3.10.12 | packaged by Anaconda, Inc. | (main, Sep 11 13:15:57) [MSC v.1916 64 bit (AMD64)]
Type "copyright", "credits" or "license()" for more information.

IPython 8.15.0 -- An enhanced Interactive Python.

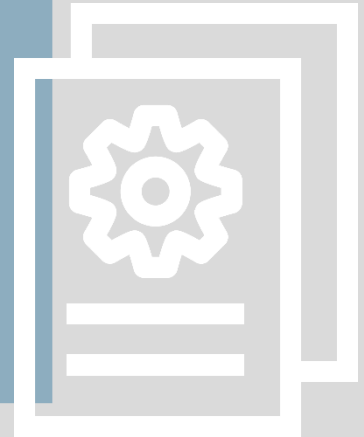
In [1]:

IPython Console History

2019 – Automation of mapping data

Lessons learned

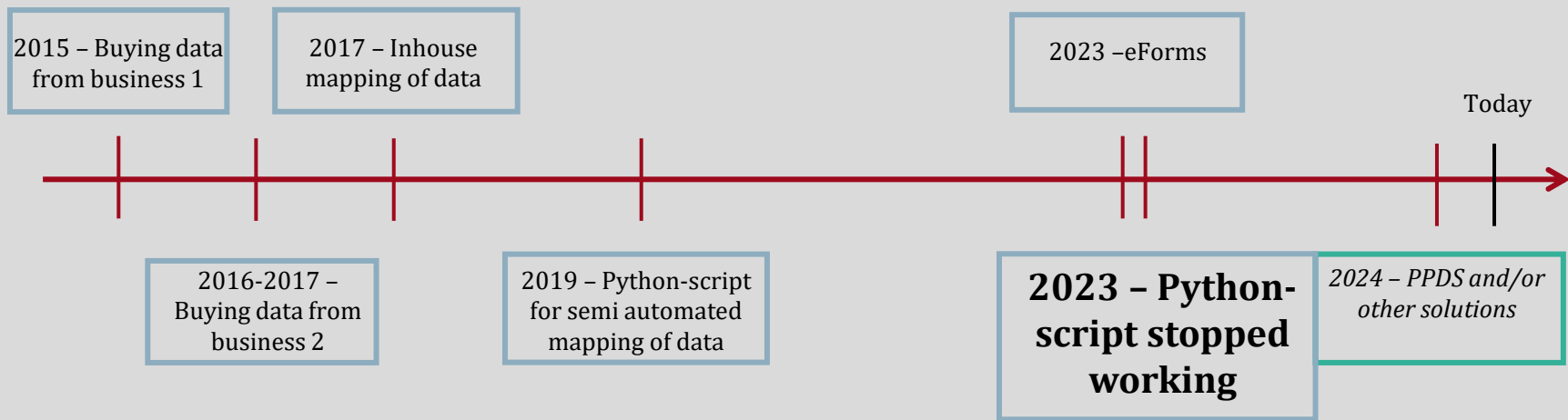
- Enormous amount of time saved
- The script optimized and streamlined mapping and quality assurance
- Additional manual quality assured an updated script and room for “weird tenders”
- Web-scraping is sensitive to change
- You need both an understanding of the filling in of notices and programming for an optimal set-up



Made available to the public in 2022

[illegible]

Timeline of procurement data



2023 – Script stopped working

Purpose and approach

- We want an automated data collection and quality assurance setup
- We are in the process of defining our new set-up

Lessons learned

- Data scientist are a limited resource
- Resources for attaining procurement data by programming must be prioritized

What have been the effects on quality?

Effects of quality assurance

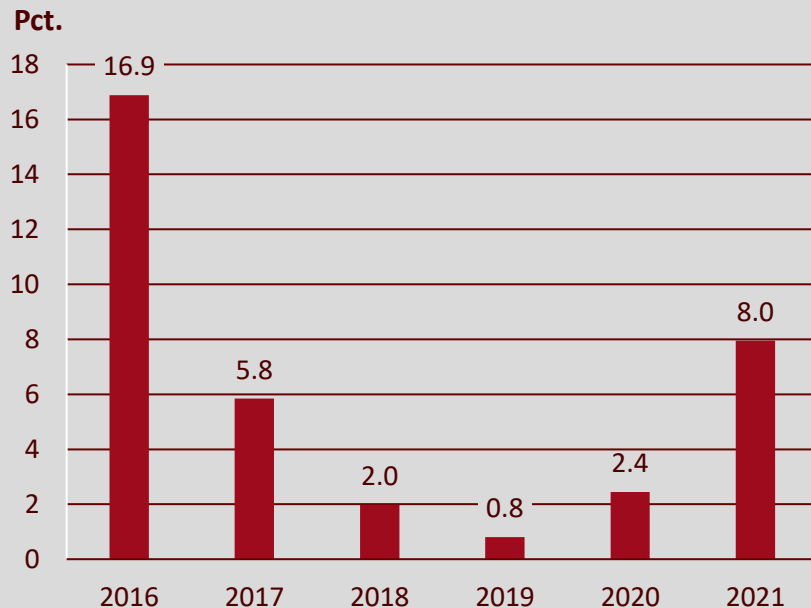
- Higher quality data
- Data has become usable for our needs
- Alterations and added information in every tender
- Measurable effects
 - Missing CAN's
 - Filled in values



Procedure for CAN-reminders

- Email to procuring authorities that have not published a CAN after 1,5 years
- High succes rate
- We found unlinked tenders and authorities published CAN's after contact
- The procedure stopped in 2021

Figure 1 - Overview of the share of missing CANs over time (2016-2021)



Fields that have been filled in

- All variables have been filled in to some extent. I have focused on a few key variables:

Total value from CAN

- 9 percent of tenders

Winning bidder

- 5 percent of tenders and lots

Business registry number

- 48 percent of tenders and lots

General effects of data mapping and quality assurance



High data quality that enables complex analyses and evaluations of our procurement law and practices



Semi-live data that can be used to inform the political level and provide legal experts with insights into practices



Awareness of data capabilities and limitations



A more complete overview of Danish tenders through increased number of CANs



Engagement and understanding amongst actual users of the notice-forms by contact



DANISH COMPETITION AND CONSUMER AUTHORITY