# The DECRYPT Project

Munich, November 1st 2019



Dr. Nils Kopal
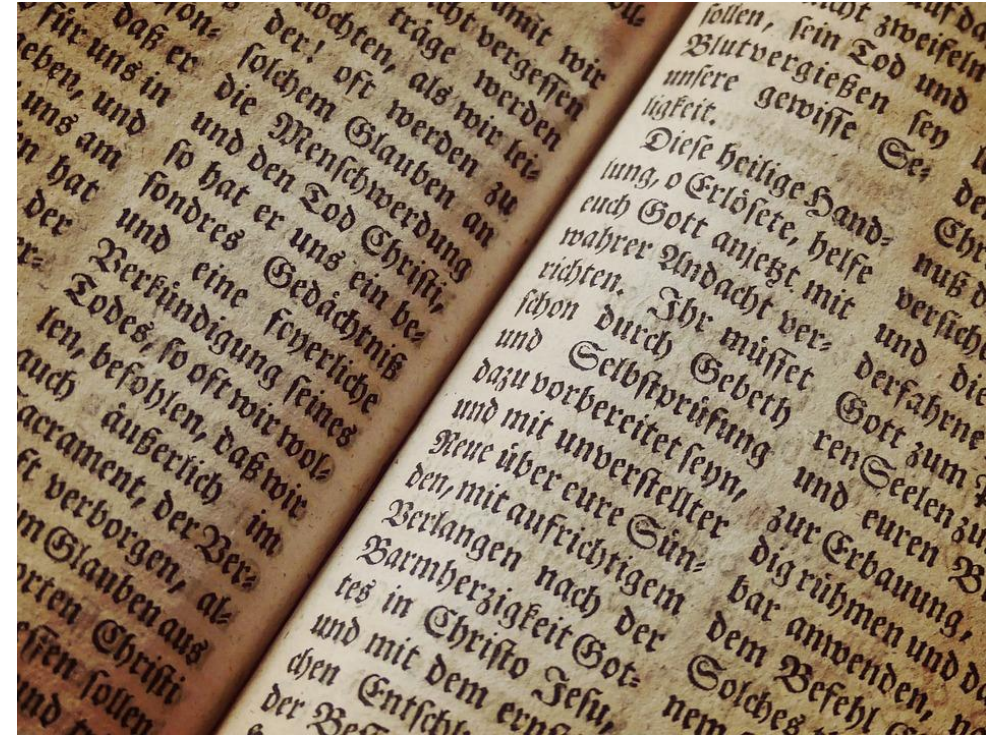
*CrypTool 2 Project / DECRYPT Project*

*kopal@cryptool.org*

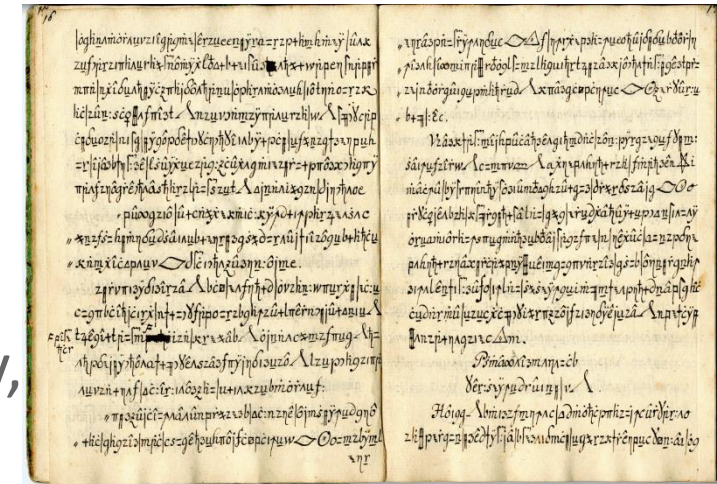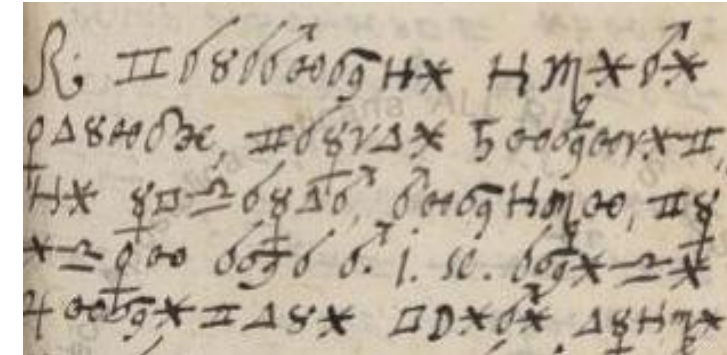*CrypTool Meeting 20+ Years*

# Content

# 1. Introduction

- **Thousands of encrypted manuscripts can be found in archives**
  - Diplomatic/military correspondence
  - Intelligence reports
  - Scientific writings
  - Private letters and diaries
  - Manuscripts related to secret societies and magic
- **Mostly not yet available for historical research**
- **Many researchers are working on these, but**
  - uncoordinated
  - in various scientific areas: history, linguistics, philology, computer science, and computational linguistics

# 1. Introduction

- **Not much is known in detail about encrypted documents throughout the centuries**
- **Large-scale studies were not possible due lack of**
    1. **infrastructural** resources
    2. **tools** for historical cryptology
- **Available tools are mostly unsuitable to these documents, because of**
    1. (old) **hand-writing**
    2. **mistakes/errors**
    3. **non-standardized languages**
    4. **no digitalization** available (**transcriptions**)
    5. often **mixtures** between **cleartext/ciphertext** and **different languages**

# 1. Introduction

- **Planned solution: DECRYPT project**

- **Project goals:**
  1. Build **infrastructural support** for historical cryptology

  2. **Collect** various documents throughout the ages, digitize them, **store** them

  3. Release **resources and tools** to **digitize**, **process** and **decrypt** historical encrypted sources

# 2. Historical Cryptology

- „**Historical cryptology is the study of encrypted messages from our history aiming at their decryption and contextualization**"

- **Dimensions**

    – **Mathematical/computerized** (cryptanalysis)

    – **Linguistic** and other coding pattern

    – **Historical** context

# 3. The DECRYPT Project

- **"DECRYPT focuses on the refinement and development of the tools involved in the automatic processing of encrypted historical sources."**

- **Supported by the Swedish Research Council, grant 2018-06074**
  - Funding for **2+4 years, started in 2019**
  - **Interdisciplinary** (linguists, historians, philologists, computer scientists, and cryptanalysts)
  - Involved universities from Sweden, Hungary, Spain, Germany (see next slide for details)

- **3 steps involved in our project:**
  1. **Data collection** and **digitization**
  2. **Analysis** needed prior to decryption
  3. **Decryption** and **cryptanalysis**

# 3. The DECRYPT Project

Project Members

- **Uppsala University, Sweden** (language analysis)
  - Beata Megyesi (project leader)
- **University of Gothenburg, Sweden** (language analysis)
  - Michelle Waldispühl
- **Computer Vision Center, Universitat Autònoma de Barcelona, Spain** (image processing)
  - Alicia Fornés
- **Budapesti Müszaki és Gazdaságtudományi Egyetem, Hungary** (history)
  - Benedek Láng
- **University of Siegen, Germany** (cryptanalysis)
  - Bernhard Esslinger
- **Universität der Bundeswehr München, Germany** (cryptanalysis)
  - Arno Wacker

- **Archive work**

  –Search for **crypto-related documents**

  –Scan/photograph

  –Cataloge (in **DECODE** database)



- **Researchers: historians**

Step 2: Analysis needed prior to decryption

- **Determine "scenario"**            **(→ attack type):**
  1. Found ciphertext            (**Ciphertext-only**)
  2. Found ciphertext and plaintext     (**Known-plaintext**)
  3. Found ciphertext and key        (**Decryption**)
  4. Found plaintext and key         (**Encryption**)


- **Researchers: cryptanalysts**

# 3. The DECRYPT Project

Step 3: Decryption including cryptanalysis

- **Part 1: (transcription)**
  - Digitization and pre-processing of the historical source resulting in **images**
  - (Semi-)automatic **transcription** of images
  - **Researchers:** image processing experts
- **Part 2a: (analysis)**
  - (Historic) **language models**
  - **Researchers**: computer linguistics, philologists
- **Part 2b: (cryptanalysis)**
  - Breaking of ciphers using e.g. **heuristics**
  - **Researchers:** cryptanalysts
- **Part 3: (historical analysis of plaintexts)**
  - Analysis of new **(historical) findings** concerning plaintexts and methods
  - **Researchers:** historians

# 4. Resources and Tools

## 1. DECODE database
– Collection of **ciphertexts, keys, etc.**

## 2. Historical corpora (HistCorp) and language models
– Collection of **historical original text corpora** in 14 European languages

## 3. Tools
– **Web service**: for transcription and „easy" parts of cryptanalysis

– **CrypTool 2**: tool for supporting „difficult" parts of cryptanalysis

– **Console applications** as prototypes → will be migrated with „nice" UI into CrypTool 2

# 5. Some Exemplary Results of the Project

- **Huge collection of historic documents in the DECODE database**
  - **~1.000 records**

  - **33% are original keys**

  - **634 cipher texts (205 decrypted, 232 transcribed)**

  - Many "one pagers"; longest document 410 pages

# 5. Some Exemplary Results of the Project

- **Decipherment of collections of historic Vatican ciphers**

| Solved from ciphertext-only | Reconstructed from plaintexts | Key found in Meister based on homophone analysis | Solved independently by N. Biermann and T. Bosbach | Unsolved | **Total** |
|---|---|---|---|---|---|
| 5 | 10 | 1 | 3 | 2 | **21** |

- **19 of 21** collections are **solved**

- **2** are **unsolved**

- A huge set of **(console) tools** developed by Lasry (and currently ported to CT2)

## Tools in CrypTool 2 – Homophonic Substitution Analyzer

## Tools in CrypTool 2 – DECODE Downloader and DECODE Viewer

## Tools in CrypTool 2 – DECODE Decipherer

# 5. Some Exemplary Results of the Project

Tools in CrypTool 2 – DECODE Symbol Heatmap

## Tools in CrypTool 2 – DECODE Key Clusterer



| DECODEClusterer | | | |
|---|---|---|---|
| **Document count:** 121 | | **Cluster count:** 49 | |

| Name | Document count | Symbol count | Different Symbols | Cluster info |
|---|---|---|---|---|
| undefined | 1 | 3 | 3 | n=33, /=33, a=33 |
| Segr. di Stato/Spagna 6_I/1/ | 1 | 454 | 10 | 0=24, 4=15, 6=12, 2=11, 5=9, 1=8, 7=8, 9=7, 8=3, 3=3 |
| Segr. di Stato Spagna 1\4 | 6 | 7344 | 10 | 0=21, 2=13, 6=13, 4=13, 5=11, 7=9, 1=9, 9=5, 8=3, 3=3 |
| Segr. di Stato Francia 346/5 | 2 | 6120 | 10 | 1=13, 3=13, 4=11, 0=11, 7=10, 8=10, 9=10, 5=8, 6=7, 2=7 |
| Segr. Stato Francia 64/1 | 3 | 61107 | 10 | 1=18, 4=14, 3=13, 0=11, 8=10, 9=9, 5=8, 7=7, 6=6, 2=5 |
| Segr. di Spagna 1\2 | 1 | 855 | 11 | 0=21, 2=14, 4=12, 6=11, 5=11, 1=9, 7=9, 9=5, 3=3, 8=2, 0^.=1 |
| Segr. Stato Spagna 364D/15 | 14 | 34070 | 10 | 8=13, 1=12, 5=12, 2=12, 7=11, 3=11, 0=10, 9=9, 6=6, 4=5 |
| Segr. Stato Spagna 1\1 | 1 | 967 | 14 | 0=23, 6=11, 4=11, 2=11, 7=9, 5=8, 1=6, 9=4, 3=4, 8=4, 6__=1, 4__=1, 2__=1, 7__=1 |
| Segr. di Stato Francia 3/1/ | 1 | 1540 | 10 | 6=18, 2=16, 8=13, 1=12, 4=11, 0=8, 3=7, 5=7, 7=3, 5^.=1 |
| Segr. Stato Spagna 6_II\12 | 6 | 7169 | 10 | 0=23, 6=13, 2=12, 4=11, 5=11, 7=9, 1=8, 9=5, 8=3, 3=3 |
| Segr. di Stato/Portogallo 117/5 | 1 | 489 | 11 | 8=20, 5=16, 4=13, 1=11, 3=10, 9=7, 7=6, 0=5, 6=4, 2=3, ==3 |
| Segr. Stato Spagna 6_II\3 | 1 | 387 | 10 | 0=21, 6=17, 2=11, 4=11, 7=10, 5=10, 1=7, 9=5, 8=3, 3=3 |
| Segr. Stato Francia 3\3\ | 1 | 377 | 18 | 2=17, 6=16, 4=10, 8=10, 1=8, 0=8, 5=6, 3=5, 7=3, 1__=2, 2__=2, 6__=2, 8__=2, 3__=2, 0__=1, |
| Segr. Stato Spagna 423\6 | 2 | 3829 | 10 | 8=13, 0=13, 3=11, 2=11, 6=10, 4=10, 7=10, 5=9, 9=9, 1=5 |
| Segr. Stato Francia 18-Split/3/ | 1 | 6374 | 10 | 1=23, 0=14, 5=13, 2=11, 7=8, 8=8, 3=7, 9=6, 4=5, 6=4 |
| Segr. Stato Francia 22/11 | 28 | 177861 | 10 | 1=23, 0=14, 5=14, 2=14, 7=10, 3=7, 9=6, 4=5, 6=5, 8=1 |
| Segr. di Stato/Portogallo 117/2 | 1 | 858 | 10 | 8=20, 5=19, 3=12, 4=10, 7=9, 1=8, 9=7, 0=7, 6=3, 2=3 |
| Segr. Stato Francia 171 | 1 | 239 | 14 | 9=19, 3=14, 1=14, 1^.=10, 0^.=8, 5=7, 4^.=6, 6^.=5, 2^.=4, 7=4, 3^.=3, 5^.=3, 6=1, 0=1 |
| Segr. Stato, Francia 41 | 1 | 1867 | 18 | 5=12, 2=11, 1=11, 6=11, 7=9, 4=7, 3=7, 0=5, 8=4, 5_.=3, 6_.=3, 4_.=3, 8_.=3, 9=2, 7^.=1, 7_. |
| Segr. di Stato Spagna 423/1 | 6 | 12797 | 10 | 6=12, 4=11, 8=11, 7=11, 2=10, 0=10, 5=10, 9=10, 3=10, 1=5 |
| Segr. di Stato Francia 7/1/ | 1 | 2917 | 27 | 2^.=13, 5=6, 1=6, i=5, a=5, 2=5, 8=4, 9=4, e=4, 3=3, "=3, r=3, o=3, 6=3, l=3, n=3, -=3, t=2, |
| Segr. Stato Francia 22/32 | 1 | 785 | 10 | 1=15, 3=13, 4=11, 9=10, 2=10, 8=9, 6=9, 5=9, 7=8, 0=6 |
| Segr. Stato Spagna 6_II\2\ | 1 | 363 | 10 | 0=24, 5=13, 6=13, 2=12, 4=12, 1=9, 7=5, 9=4, 3=4, 8=2 |
| Segr. Stato Spagna 364C/18 | 9 | 38595 | 10 | 8=13, 1=13, 5=13, 2=12, 7=12, 0=11, 9=9, 6=6, 4=6, 3=5 |

# 5. Some Exemplary Results of the Project
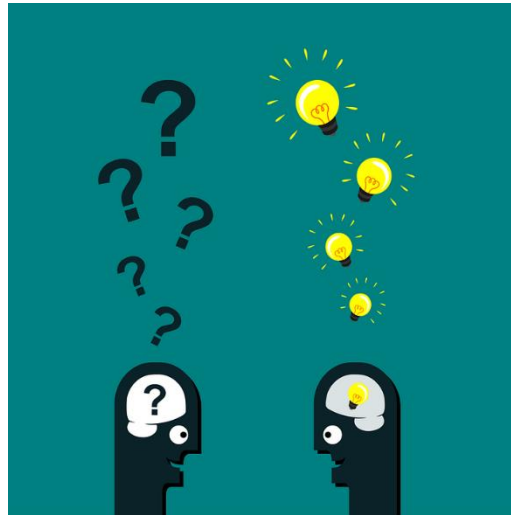
Tools in CrypTool 2 – DECODE Parser Tester

# 6. Conclusion

- **DECRYPT collects historic ciphers and develops tools and resources for transcription and cryptanalysis**

- **Project limits & challenges**
  - Will **NOT** be able to **solve every cipher automatically**, e.g. non-deterministic ciphers
  - Sometimes, **manual transcriptions are more feasible** than automatic due to errors
  - "**Unknown**" cipher (types) will also **be hard** to be solved automatically
  - Tools often will **support the cryptanalyst** and can not replace him ☺

- **Benefits**
  - Huge set of **original historic crypto material**
  - Many **helpful tools**
  - New insights in **early-modern cryptology**
  - New insights in our **"hidden history"**

# Questions and discussion

**Thank you very much for your attention!**



**Do you have questions?**