



UPPSALA
UNIVERSITET

Uncovering secrets from the past:

Research of historical cipher keys within the DECRYPT project



Crina Tudor
Department of Linguistics and Philology
Uppsala University

January 31st, 2022



Historical ciphers

- Thousands of ciphers are buried in archives
 - Diplomatic correspondence
 - Intelligent reports
 - Docs from secret societies
 - Private letters and diaries
- Not indexed, not marked up as such

- ✓ Need to bring to light the content for historical contextualization
- ✓ Require research infrastructure for historical cryptology
- ✓ Need of large-scale, systematic studies





UPPSALA
UNIVERSITET

The DECRYPT project

The DECRYPT Team

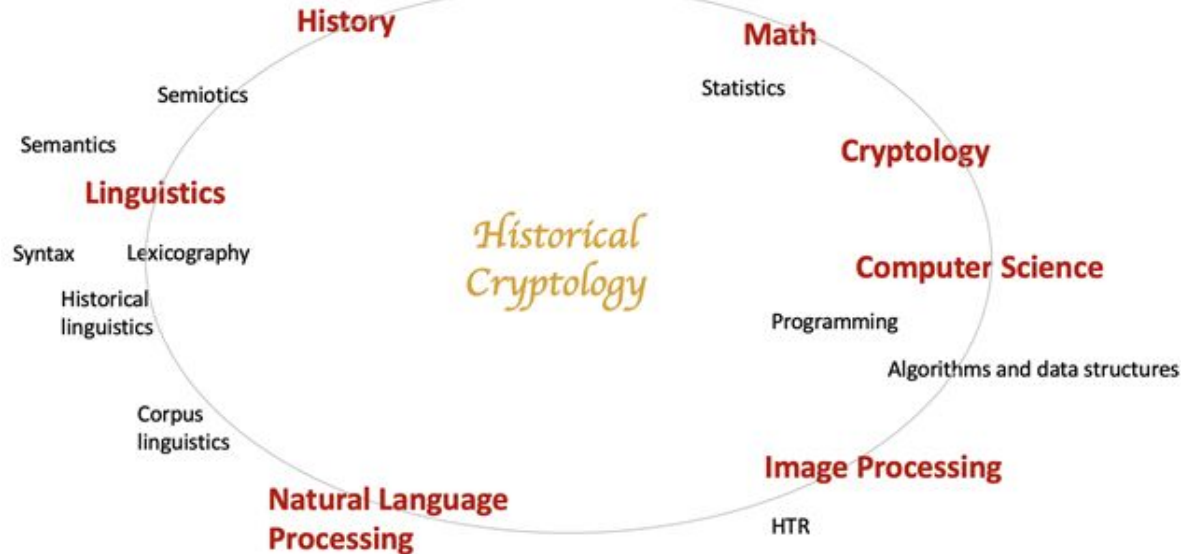
<u>Benedek Láng</u> history	<u>Michelle Waldispöhl</u> historical linguistics	<u>Eva Pettersson</u> NLP	<u>Beáta Megyesi</u> NLP, PI	<u>Alicia Fornés</u> HTR	<u>Bernhard Esslinger</u> crypto	<u>Arno Wacker</u> crypto
<u>Anna Lehofer</u> history	<u>Karl de Leeuw</u> history	<u>Crina Tudor</u> NLP	<u>Mihály Héder</u> sysdev	<u>Nils Kopál</u> crypto	<u>George László</u> crypto	<u>Vasily Mikhalev</u> crypto
<u>Mohamed Ali Souibgui</u> HTR	<u>Jialou Chen</u> HTR	<u>Arnau Baró</u> HTR	<u>Yousri Kessentini</u> HTR	<u>Ferenc Sziéeti</u> Transcription Tool	<u>Giacomo Magnifico</u> Transcription Evaluation	historians, librarians, students, ...

Computer Vision Center at Universitat Autònoma de Barcelona.



Project aim

Establish a new cross-disciplinary scientific field of historical cryptology to decrypt and contextualize historical encrypted sources.





The DECRYPT Portal

<https://de-crypt.org/index.php>

Resources:

- [The DECODE database](#)
- [HistCorp](#)

Tools:

- [CryptTool2](#) - cipher breaking
- [Transcript](#) - interactive transcription tool
- [Decoder](#) - ciphertext-key mapping
- [Anacode](#) - ciphertext analysis
- [Anakey](#) - cipher key analysis



Anakey - purpose

- To provide a reliable transcription scheme for the transcription of historical keys
- To build a method for automatically identifying different types of keys
- To provide a statistical analysis for individual keys



Background

- No previous study conducted on the structure of historical keys
- Most of the specialized literature focuses on ciphers rather than keys
- The few studies that discuss transcription methods for keys or ciphers tend to focus on individual instances and are not conducted on a large scale
- A completely automatic and accurate transcription is currently not achievable through OCR alone



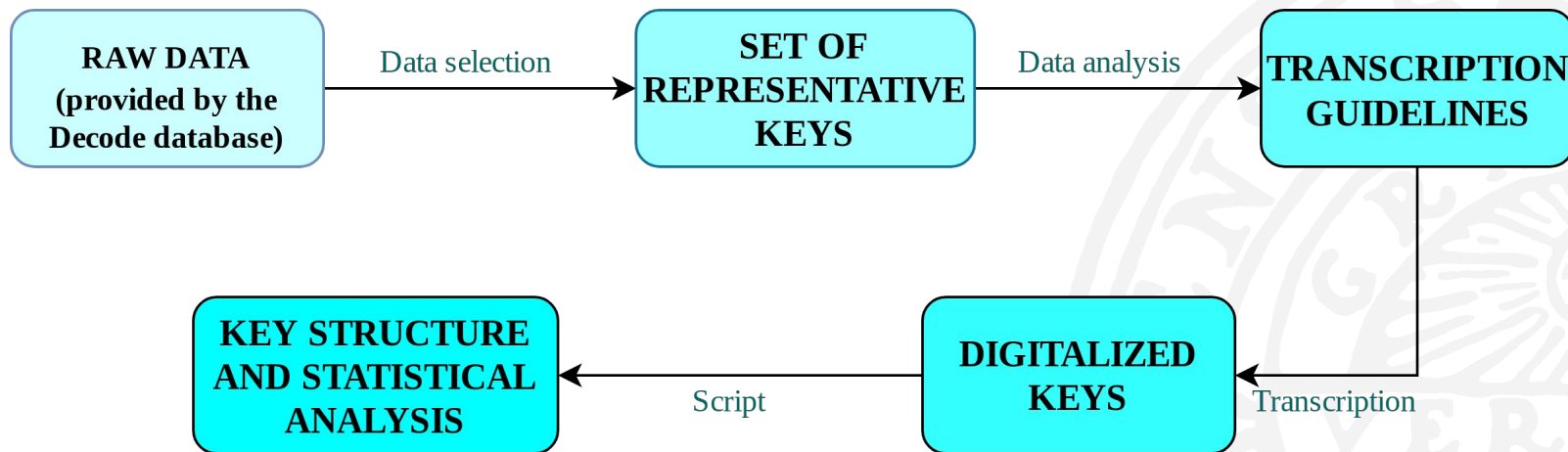
UPPSALA
UNIVERSITET

Key excerpt





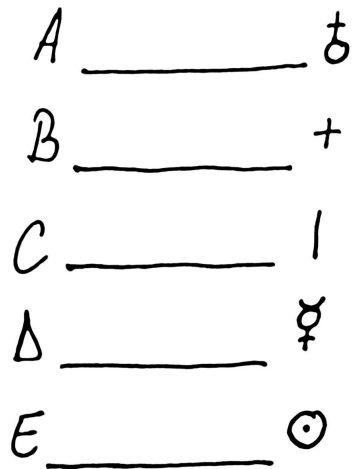
Method





Transcription Conventions

- Each transcription file is preceded by metadata
- We focus on methods for simple, homophonic and polyphonic substitution
- We focus on the non-ASCII characters
- We differentiate between 3 major symbol sets:
Latin alphabet, digits, and graphic signs



A _____ †
B _____ +
C _____ |
Δ _____ ♀
E _____ ⊙



Polyphonic substitution

<i>As</i>	<i>tn</i>	<i>mo</i>	<i>im</i>	<i>lu</i>	<i>ce</i>	<i>pd</i>	<i>bz</i>	<i>fg</i>
3	6	5	8	9	7	0	02	04

Homophonic substitution

Key excerpt



UPPSALA
UNIVERSITET

Key excerpt

#KEY: original
#CATALOG NAME: TNA_SP106/2_ElizabethI_f58(0069)
#IMAGE NAME: 3391.jpg
#LANGUAGE: FR EN
#TRANSCRIBER NAME: CT
#DATE OF TRANSCRIPTION: 10.04.2019
#TRANSCRIPTION TIME: 2h
#STATUS:complete

m - a
n - b
o - c
p - d
q - e
r - f
s - g
t - h
u - i
w - k
x - l
y - m
z - n
a - o
b - p
c - q
d - r
e - s
f - t
g - u
h - w
i - x
k - y
l - z





Automatic Key Structure Extraction

- We build an automatic method for extracting key statistics
- Some of the information we can extract using our script includes:
 - type of symbols used for encryption
 - code structure (ngraphs)
 - unknown symbols
 - plaintext structure (ngrams)
 - code distribution



Error catching

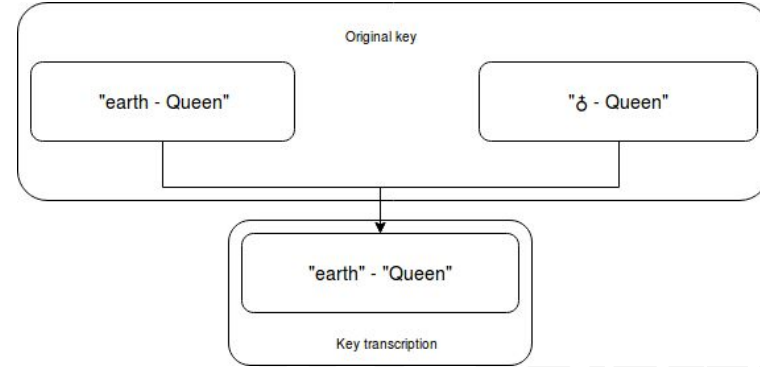
Our script can catch the following types of errors:

- metadata error LANGUAGE: FR EN > #LANGUAGE: FR EN
- delimitation error 89 a > 89 - a
- spacing error 89 - a > 89 - a
- other



Future Work

- Our method provides a solid basis for further studies into the structure of keys
- We can still benefit from certain improvements:
 - eliminating the ambiguity when it comes to graphic signs
 - the way we handle sloppy keys
 - expanding our approach to other encryption methods



Duca di Fiorenza	84	come	50
Duca di Ferrara	94	quando	60
Duca d'Urbino	22	qua	70
Duca Ottavio	32	que	80
Cote Federico Bonomes	52	qui	90
Conte Amiballe d'emp	62	queste	82
Duca d'Alma	72	quelle	84
Duca di sessa	92	Non	86
Marchese di Mondragon	24	Lettere	88
Conti di Fera	34	Nagote	82



Conclusions

- We provide a dependable transcription standard for historical keys
- We build a method for automatic key structure extraction
- We build a reliable basis for further large-scale comparative studies



Large-scale study

Megyesi B, Tudor C, Láng B, Lehofer A. *Key Design in the Early Modern Era in Europe*.
In: Proceedings of the 4th International Conference on Historical Cryptology (HistoCrypt 2021)

- *“At first, the substitution symbols were neither letters or numbers but fanciful signs like % or . . . But nobody has looked into when, in the later evolution, as nomenclators ran out of easily distinguishable symbols and began using numbers, the cipher secretaries began forming two-part nomenclators. This research requires merely examining the many nomenclators in the archives of Italy and France and timing and quantifying the change. I suppose it will be tough, living in Europe for a year and having an aperitif after a day examining antique manuscripts. But somebody should do it!” (Kahn, 2008:58)*
- This is exactly what we are up to....



Research questions

- What types of keys were used in Europe between the 15th and 18th centuries? What were their specific characteristics?
- What was encoded and how?
- How did encryption evolve over time?
- Can we apply simple statistical methods to large-scale analysis of transcribed historical keys?



Goal

- provide insight into the evolution of encryption
 - pilot: original keys from ca 1500-1800 in Europe
- provide a structural description of keys along with their morphological analysis and a typology
- investigate the key structure in terms of:
 - what is encoded - plaintext: languages, entity types
 - how is encoded - ciphertext: code types, symbols systems
- describe some trends of the code structure related to: time periods, geographic areas, ...



Method

- transcription guidelines for keys
 - common symbol representation across keys
- structural description of keys:
 - morphological (inner) structure
 - definition of key types
 - automatic structural description of keys
- match metadata info and structural description



Data

- 1665 keys available in the DECODE collection from the 15th-18th centuries from Austria, Belgium, France, Germany, Hungary, Italy, Spain, the Netherlands, the UK, and the Vatican
- 26% of the keys have been manually transcribed



Automatic extraction of key structure

Extracted features from original key transcriptions:

- symbols types (digits, Latin, Greek, graphic signs, ...)
- code types:
 - no. of unique code types
 - unigraph, digraph, trigraph, 4+graph
 - fixed vs variable length
- plaintext types:
 - total number of unique plaintext units
 - unigrams, bigrams, trigrams, 4+grams

Key excerpt

Key excerpt



Cipher type

Cipher type:

- defined on several levels:
alphabet and nomenclature
- simple substitution
- homophonic substitution
- polyphonic substitution

Key excerpt

Key excerpt

<i>As</i>	<i>tn</i>	<i>mo</i>	<i>im</i>	<i>lu</i>	<i>ce</i>	<i>pd</i>	<i>bz</i>	<i>fg</i>
3	6	5	8	9	7	0	02	04

Key excerpt



Process

Key excerpt



```
#KEY: original
#CATALOG NAME: ASV_ARM_XLIV_7-1
#IMAGE NAME: 1287.jpg
#LANGUAGE: IT
#TRANSCRIBER NAME: CT
#DATE OF TRANSCRIPTION: 19.04.2019
#TRANSCRIPTION TIME: 2h
#STATUS:complete

#TYPE-NC: people, geographic names, common words

3 - A|s
6 - t|r
5 - n|o
8 - i|m
9 - l|u
7 - c|e
0 - p|d
02 - b|z
04 - f|g
00 - &con

22 - Il PP. N.S.,S.S.^ta
32 - Imp. S. M^ta Ces^a
42 - Re Filippo, Re Cath^co S. di^ta Cath^ca
52 - Re di Francia, Re dn^m s. di^ta Chr^ma
62 - Re di Portogallo
72 - Re di Polonia
82 - Re di Bohemia, S.ses^ta|queste|Anist
92 - Il Principe di Spagna
24 - Mons' di Vandomo
34 - Mad. la Regente, Regina di Francia
44 - Principessa di Portogallo
54 - Sig. Venetiani
64 - Sig. Suizze
74 - Duca di Savoia
84 - Duca di Fiarenza|quelle
94 - Duca di Ferrara
22^ - Duca d'Urbino
32^ - Duca Ottavio
52^ - Co'te Federico Borromeo
```



```
Metadata
#CATALOG NAME: ASV_ARM_XLIV_7-1
#LANGUAGE: IT
#STATUS:complete
#TRANSCRIPTION TIME: 2h
#TYPE-NC: people, geographic names, common words

Cipher symbols:dgilts, graphtc signs

Total number of unique ciphertext symbols:78
  unigraphs:12
    out of which dgilts:7
  digraphs:66
    out of which dgilts:66
  trigraphs:0
  4graphs:0

Total number of ciphertext symbols matched:78
No new ciphertext symbols were found.

Total number of unique plaintext units:94
  out of which unigrams:18
  out of which bigrams:0
  out of which trigrams:7
  out of which 4grams:69
  out of which nulls:0
  out of which empty:0

Code distribution

Cipher type:polyphonic substitution

Code type:variable length

Number of codes encoding plaintext
  unigrams:9
  bigrams:0
  trigrams:6
  4grams:64
  nulls:0
  empty:0

Distribution according to plaintext type (ciphertext:plaintext)
  1. Alphabet
  The alphabet has a uniform 1:1 distribution.
  2. Nomenclature
    1:1 65
    1:2 4
    1:3 1
    1:4+ 0
```




CSV output

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	AA	AB
Catalog name	Origin city	Origin country	Date	Language	Cipher symbols		Total # of unique ciphertext symbols	unigrams	unigram digits	digrams	digram digits	trigrams	trigram digits	4-grams	4-gram digits	Total number of unique plaintext units	unigrams	bigrams	trigrams	4-grams	n-grams	entropy	cancellation		Cipher type	Code type	
1	N/A	N/A	1703-01-01 - 1711-12-31	LA, EN	digits		298	0	0	58	58	238	238	1	1		254	71	93	8	127	6	0	0		homophonic substitution, polyphonic substitution, simple substitution	variable
2	NAH_015_CAPS_C_FASC_43_3	N/A	1703-01-01 - 1711-12-31	LA	digits		301	9	9	90	90	201	201	0	0		295	33	111	16	134	7	0	0		homophonic substitution, simple substitution	variable
3	NAH_015_CAPS_C_FASC_43_2	N/A	1703-01-01 - 1711-12-31	LA	digits		94	9	9	75	75	10	10	0	0		45	96	1	1	17	4	0	0		homophonic substitution, polyphonic substitution, simple substitution	variable
4	RAJ_ARCH_Rakoczi_C64_462_29_3	N/A	1703-01-01 - 1711-12-31	LA	digits		12	0	0	11	11	1	1	0	0		3	8	4	0	0	0	0	0		homophonic substitution	variable
5	RAJ_ARCH_Rakoczi_C64_462_29_9	N/A	1703-01-01 - 1711-12-31	LA	digits		169	9	9	88	88	61	61	0	0		123	50	3	0	98	8	0	0		homophonic substitution	variable
6	THA_SPI064_Linnel119-200005-0004	N/A	1614-01-01 - 1617-12-31	EN	digits		231	0	0	70	70	161	161	0	0		182	69	118	23	22	0	0	0		homophonic substitution	variable
7	RAJ_ARCH_Rakoczi_C64_462_29_4	N/A	1703-01-01 - 1711-12-31	LA	digits		300	21	1	74	72	5	5	0	0		100	24	47	4	24	1	0	0		homophonic substitution, simple substitution	variable
8	Flemming in Vienna Saxony 1761_HSID_10024_Loc_0823611_Bi_24	Denmark	Saxony	1761-01-01 - 1761-12-31	FR	digits, Latin alphabet	266	0	0	66	66	200	200	0	0		217	78	85	6	100	0	0	0		homophonic substitution, simple substitution	variable
9	RAJ_ARCH_Rakoczi_C64_462_29_5	N/A	1703-01-01 - 1711-12-31	LA	digits		666	0	0	0	0	666	666	0	0		668	0	575	2	91	0	0	0		simple substitution	fixed
10	THA_SPI067_Anne_(0087-0089)-1	N/A	1702-03-08 - 1817-06-20	FR	digits		678	0	0	0	0	678	678	0	0		676	0	575	0	99	0	0	0		simple substitution	fixed
11	THA_SPI067_Anne_(0087-0089)-2	N/A	1702-03-08 - 1817-06-20	FR	digits		356	7	7	90	90	199	114	58	0		293	2	56	100	178	20	0	0		homophonic substitution	variable
12	THA_SPI068_C2095611_0163-0164	N/A	1672-01-01 - 1672-12-31	LA	digits		437	0	0	56	56	381	381	0	0		418	43	109	19	266	1	0	0		homophonic substitution, simple substitution	variable
13	THA_SPI068_C2095611_0023-0024	N/A	1660-05-29 - 1665-00-06	EN	digits		198	0	0	90	90	105	105	0	0		196	40	123	7	26	0	0	0		homophonic substitution, polyphonic substitution	variable
14	RAJ_ARCH_Rakoczi_C64_462_29_29	N/A	1703-01-01 - 1711-12-31	LA	digits		431	9	9	88	88	334	334	0	0		382	58	172	4	182	17	0	0		homophonic substitution	variable
15	THA_SPI065_C1095611_0008-0009	N/A	1630-01-01 - 1650-12-31	EN	digits		30	5	5	25	25	0	0	0	0		32	28	1	1	2	0	0	0		homophonic substitution, simple substitution	variable
16	Saxony 1761_HSID_10024_Loc_0823611_I_29-30	Denmark	Saxony	1761-01-01 - 1761-12-31	FR	digits	114	10	10	96	96	8	8	0	0		130	24	46	6	41	2	0	0		homophonic substitution, polyphonic substitution	variable
17	Copenhagen-Saxony-Hamburg 1761_HSID_10024_Loc_0823611_I_28-29	Denmark	Saxony	1761-01-01 - 1761-12-31	FR	digits	107	34	9	73	71	0	0	0	0		108	23	52	4	28	0	0	0		polyphonic substitution	variable
18	Cakenberg-Saxony 1761_HSID_10024_Loc_0823611_Bi_3	N-CITY: Dresden	COUNTRY: Sax	1761-01-01 - 1761-12-31	FR	digits, Latin alphabet																					
19																											
20																											
21																											
22																											
23																											
24																											
25																											
26																											
27																											
28																											
29																											





Trends

Goal:

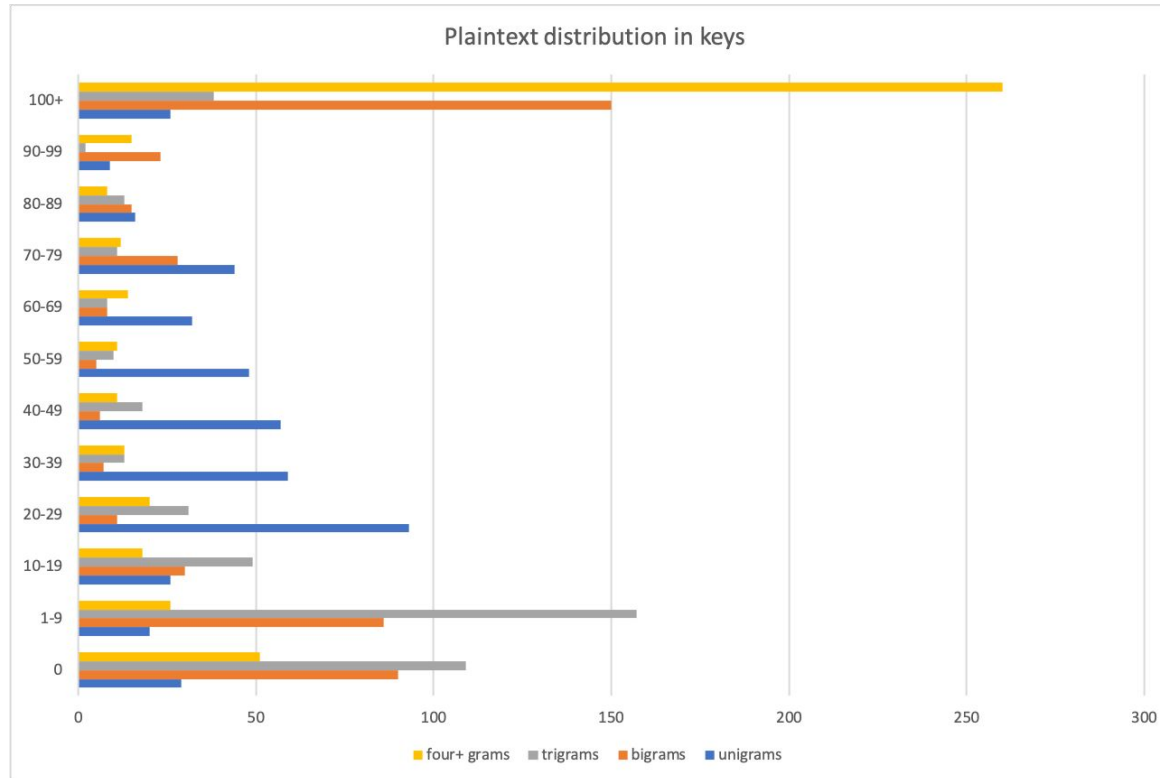
- investigate the trends throughout the 15th-18th centuries
- what have been chosen to be encoded and how

Data:

- 450 keys: automatically extracted from transcriptions
- 250 keys: manually extracted structural information without any transcriptions originating from the 15th and 16th centuries.
- 700 keys in total

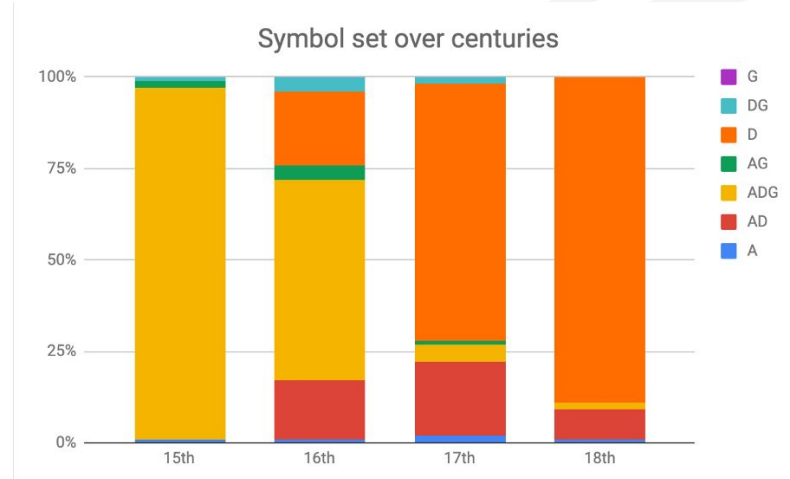
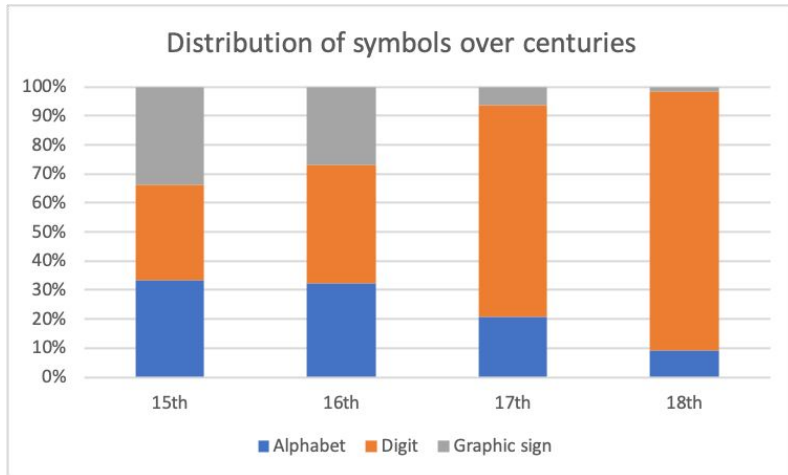


Plaintext





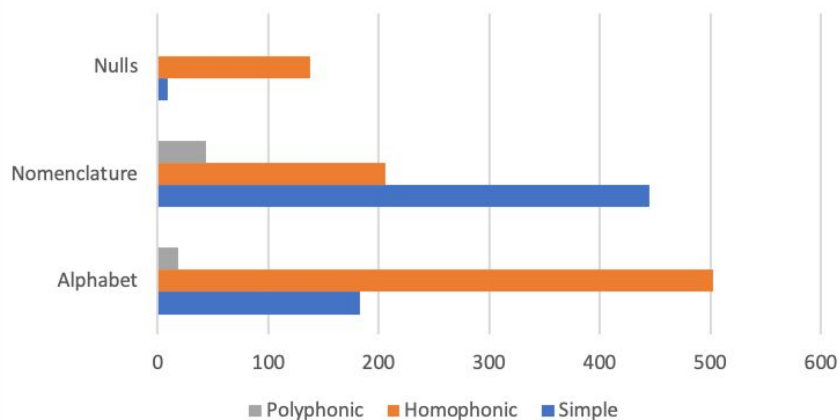
Symbol set



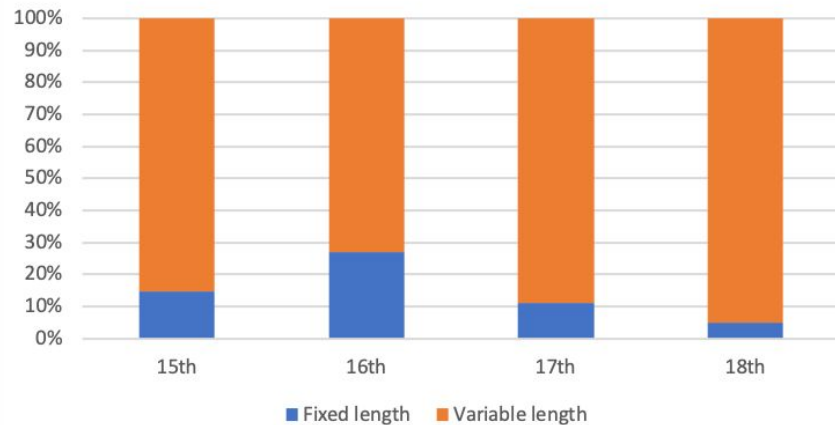


Codes

Code types in keys



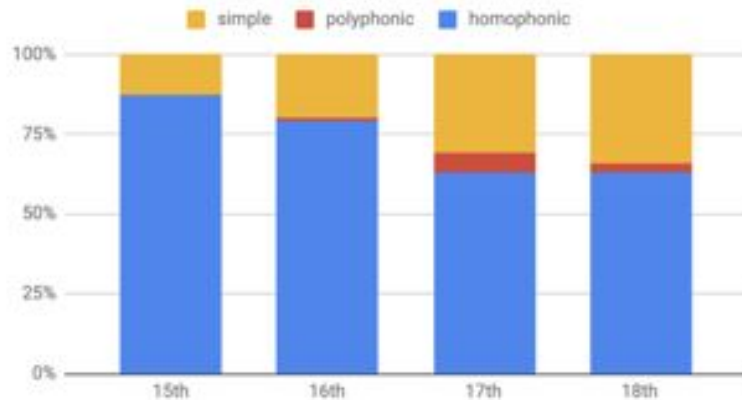
Code length over centuries



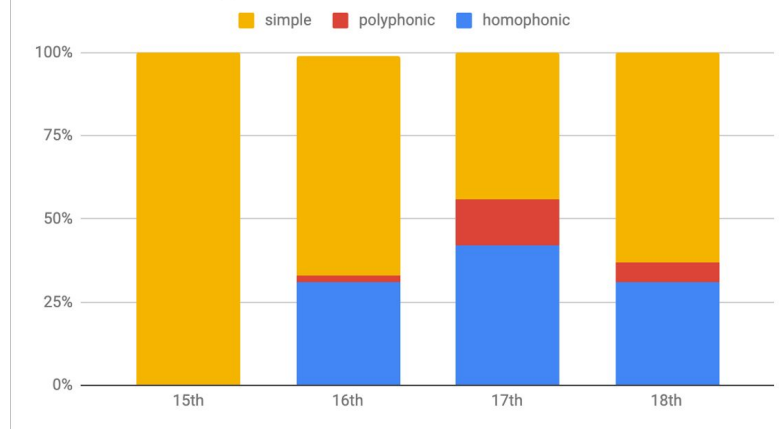


Codes

Code types for alphabetical signs over centuries

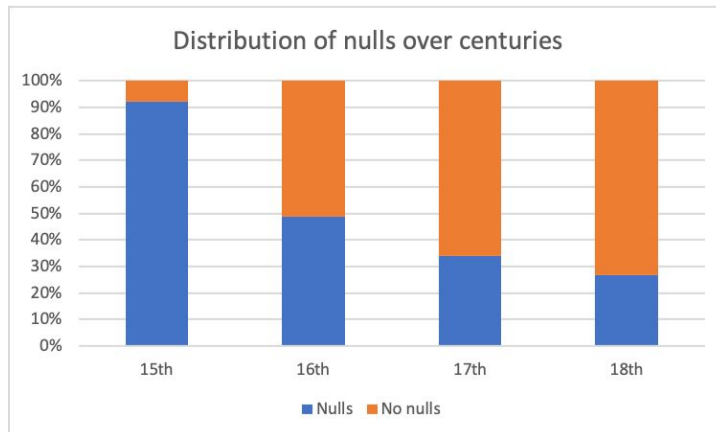


Code types for nomenclatures over centuries





Nulls





Conclusion

- We investigated 700 cipher keys from the 15th to the 18th centuries
- We described the keys' internal structure and their morphology:
 - what have been chosen to be encoded and how
 - the type of the symbol set and the code structures used, and the changes and trends of each century.
- Keys evolved over time and their structure changed



Conclusion

- Codes with various symbols including alphabets, digits, and graphic signs were dominating in the 15th century, digits only became more frequent and became the standard in the 18th century.
- The codes varied in length for alphabetical signs and nomenclatures throughout all centuries and codes with fixed length seemed to be most popular in the 16th century.
- Coding alphabetical signs was mostly homophonic, but simple substitution of letters became more frequent as the length of the nomenclatures increased over time.
- Nomenclatures were mostly encoded as simple substitution.
- Nulls have been frequently used in the 15th century and decreased significantly over time.
- Cancellation as phenomenon became "popular" in the 18th century.



Future/upcoming work

- include more data
- more precise metadata with location (GIS) (and person)!
 - correlations between features
 - add automatic key complexity (cont.)
 - JSON representation of automatic structural description
 - include the automatic structural description into the decrypt pipe **ANAKEY**



UPPSALA
UNIVERSITET

Thank you!

