

The Copiale Cipher*

Kevin Knight

USC Information Sciences Institute
4676 Admiralty Way
Marina del Rey, CA, 90292, USA
knight@isi.edu

SDL Language Weaver, Inc.
6060 Center Drive, Suite 150
Los Angeles, CA 90045
kknight@sdl.com

Beáta Megyesi and Christiane Schaefer

Department of Linguistics and Philology
Uppsala University
Box 635, S-751 26 Uppsala, Sweden
beata.megyesi@lingfil.uu.se
christiane.schaefer@lingfil.uu.se

Abstract

The Copiale cipher is a 105-page enciphered book dated 1866. We describe the features of the book and the method by which we deciphered it.

1. Features of the Text

Figure 1 shows a portion of an enciphered book from the East Berlin Academy. The book has the following features:

- It is 105 pages long, containing about 75,000 handwritten characters.
- The handwriting is extremely neat.
- Some characters are Roman letters (such as **a** and **b**), while others are abstract symbols (such as **ϣ** and **Δ**). Roman letters appear in both uppercase and lowercase forms.
- Lines of text are both left- and right-justified.
- There are only a few author corrections.
- There is no word spacing.

There are no illustrations or chapter breaks, but the text has formatting:

- Paragraphs are indented.
- Some lines are centered.

—
**This material was presented as part of an invited talk at the 4th Workshop on Building and Using Comparable Corpora (BUCC 2011).*

- Some sections of text contain a double-quote mark (") before each line.
- Some lines end with full stop (.) or colon (:). The colon (:) is also a frequent word-internal cipher letter.
- Paragraphs and section titles always begin with Roman letters (in capitalized form).

The only non-enciphered inscriptions in the book are “Philipp 1866” and “Copiales 3”, the latter of which we used to name the cipher.

The book also contains preview fragments (“catchwords”) at the bottom of left-hand pages. Each catchword is a copy of the first few letters from the following (right-hand) page. For example, in Figure 1, the short sequence **ϣâλ** floats at the bottom of the left page, and the next page begins **ϣâλomi...** In early printing, catchwords were used to help printers validate the folding and stacking of pages.

2. Transcription

To get a machine-readable version of the text, we devised the transcription scheme in Figure 2. According to this scheme, the line

πδ|ρϣΔϣêϣcâ=ϣλûb⊕ur⊗ε

is typed as:

pi oh j v hd tri arr eh three c. ah
ni arr lam uh b lip uu r o.. zs

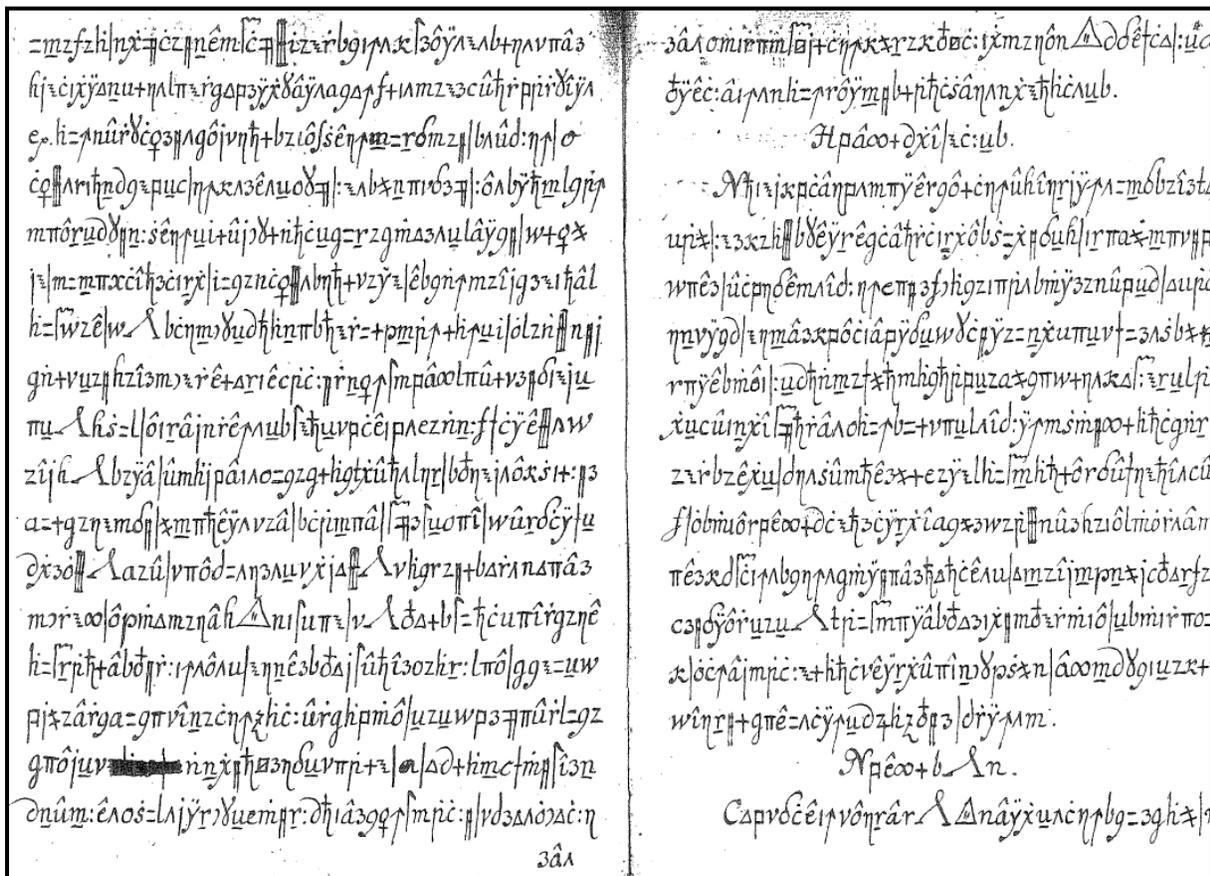


Figure 1. Two pages from the Copiale cipher.

The transcription uses easy-to-reach keyboard characters, so a transcriber can work without taking his/her eyes off the original document.

There are approximately 90 cipher letters, including 26 unaccented Roman letters, a-z. The letters c, h, m, n, p, r, s, and x have dotted forms (e.g., \dot{c}), while the letter i also has an un-dotted form. The letters m, n, r, and u have underlined variants (e.g., \underline{m}), and the vowels have circumflexed variants (e.g., \hat{a}). The plain letter y does not appear unaccented until page 30, but it appears quite frequently with an umlaut (\ddot{y}). The four Roman letters d, g, n, and z appear in both plain (\mathfrak{d} , \mathfrak{g} , \mathfrak{n} , \mathfrak{z}) and fancy forms (\mathfrak{d} , \mathfrak{g} , \mathfrak{n} , \mathfrak{z}). Capitalized Roman letters are used to start paragraphs. We transcribe these with A-Z, though we down-case them before counting frequencies (Section 3). Down-casing D, G, N, and Z is not trivial, due to the presence of both plain and fancy lowercase forms.

The non-Roman characters are an eclectic mix of symbols, including some Greek letters. Eight symbols are rendered larger than others in the text: Λ , Θ , Δ , \times , \circ , \mathfrak{A} , ∞ , and Γ .

We transcribed a total of 16 pages (10,840 letters). We carried out our analysis on those pages, after stripping catchwords and down-casing all Roman letters.

3. Letter Frequencies and Contexts

Figure 3 shows cipher letter frequencies. The distribution is far from flat, the most frequent letter being \wedge (occurring 412 times). Here are the most common cipher digraphs (letter pairs) and trigraphs, with frequencies:

\mathfrak{r} \mathfrak{h}	99	\mathfrak{r} \mathfrak{h} \wedge	47
\dot{c} :	66	\dot{c} : \underline{u}	23
\mathfrak{h} \wedge	49	η \mathfrak{r} \mathfrak{h}	22
: \underline{u}	48	\ddot{y} \mathfrak{r} \mathfrak{h}	18
\mathfrak{z} \mathfrak{h}	44	\mathfrak{h} \dot{c}	17

a	a	â	ah			δ	del
b	b					Δ	tri
c	c			ç	c.	ϝ	gam
d	d					ι	iot
e	e	ê	eh			^	lam
f	f					π	pi
g	g					ʄ	arr
h	h	ĥ	h.	h̃	hd	ʔ	bas
i	i	î	ih			ʔ	car
j	j					+	plus
k	k					†	cross
l	l					♀	fem
m	m	m̂	m.	m̃	mu	♁	mal
n	n	n̂	n.	ñ	nu	♁	ft
o	o	ô	oh	ó	o.	⊠	no
p	p	p̂	p.			♁	sqp
q	q					z	zzz
r	r	r̂	r.	r̃	ru	f	pipe
s	s	ŝ	s.			f	longs
t	t					fl	grr
u	u	û	uh	ũ	uu	fl	grl
v	v					fl	grc
w	w			Δ	tri..	ʔ	hk
x	x	x̂	x.	⊠	lip	Γ	sqi
y	y	ŷ	y..	λ	nee	:	:
z	z			⊠	o..	.	.
ð	ds	=	ni	♁	star
g	gs	κ	ki	κ	bigx		bar
z	zs	€	smil	Π	gat	3	three
η	ns	;)̂	smir	∞	toe	∞	inf

Figure 2. Transcription scheme. Columns alternate between the cipher letters and their transcriptions.

The full digraph counts reveal interesting patterns among groups of letters. For example, letters with circumflexes (â, ê, î, ô, û) have behaviors in common: all five tend to be preceded by z and π, and all five tend to be followed by 3 and j. To get a better handle on letter similarities, we automatically clustered the cipher letters based on their contexts. The result is shown in Figure 4. We did the clustering as follows. For each distinct letter x, we created a

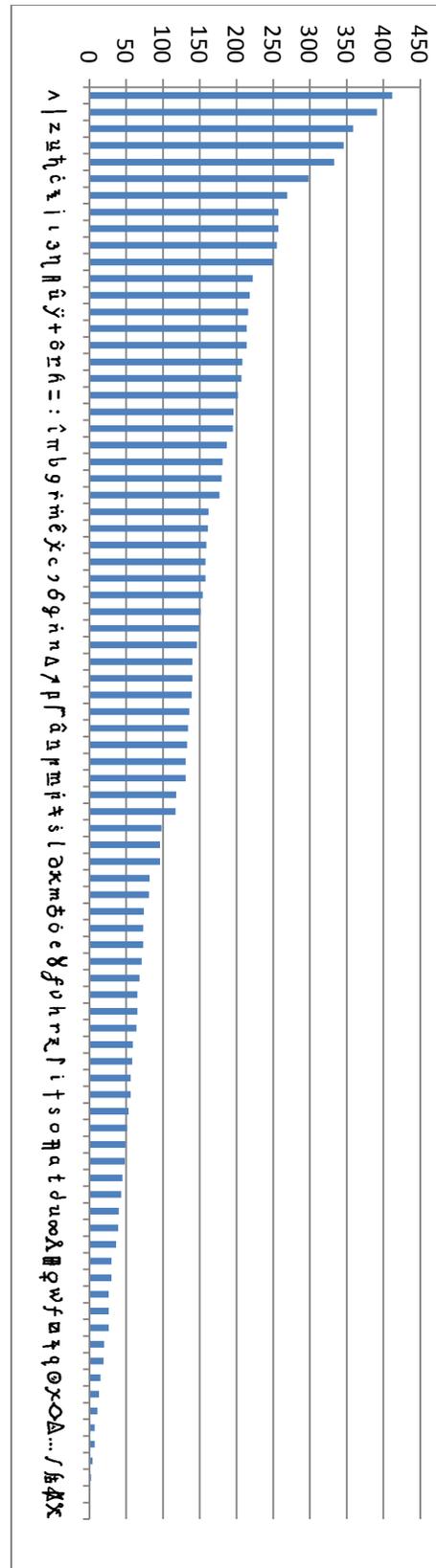


Figure 3. Cipher letter frequencies.

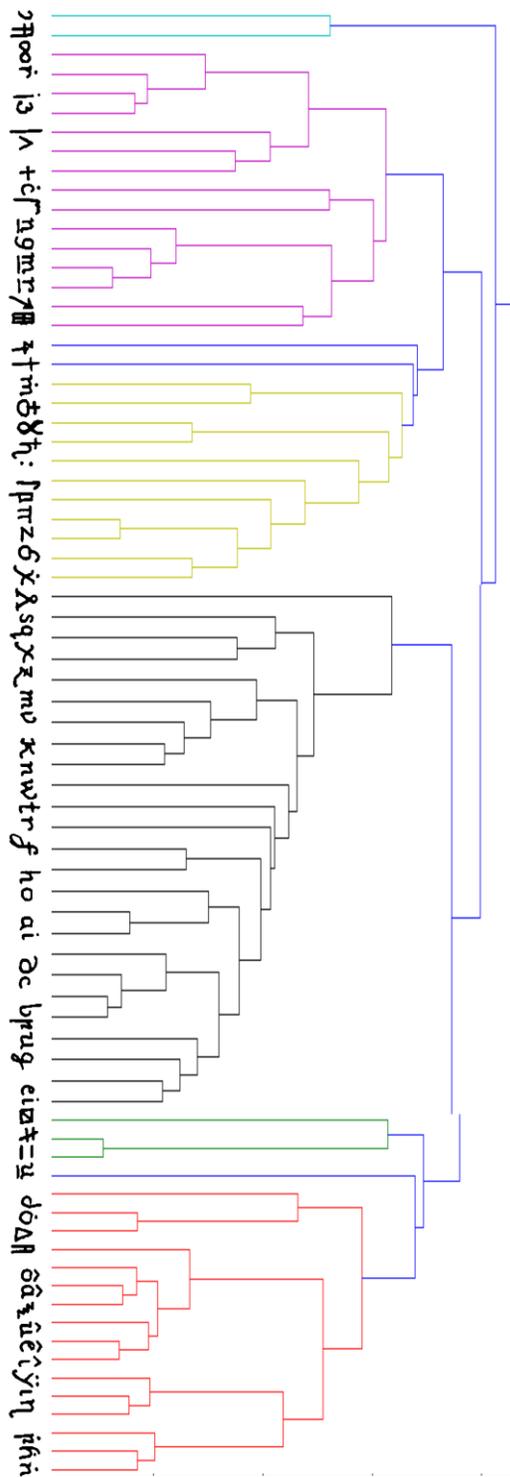


Figure 4. Automatic clustering of cipher letters based on similarity of contexts.

co-occurrence vector of length 90, to capture the distribution of letters than precede x . For example, if x is preceded 12 times by π , 0 times by \hat{u} , 4 times by \tilde{y} , 1 time by δ , etc, then its vector looks like this: $[12, 0, 4, 1, \dots]$. For the same letter x , we created another vector that captures the distribution of letters than follow x , e.g., $[0, 0, 7, 2, \dots]$. Then we concatenated the two vectors to create $v(x) = [12, 0, 4, 1, \dots, 0, 0, 7, 2, \dots]$. We deemed two letters a and b to be similar if the cosine distance between $v(a)$ and $v(b)$ is small, indicating that they appear in similar contexts. We used the Scipy software (<http://users.soe.ucsc.edu/~eads/cluster.html>) to perform and plot a clustering that incrementally merges similar letters (and groups of letters) in a bottom-up fashion.

The cluster diagram confirms that circumflexed letters (\hat{a} , \hat{e} , \hat{i} , \hat{o} , \hat{u}) behave similarly. It also shows that the unaccented Roman letters form a natural grouping, as do underlined letters. Merges that happen low in the cluster map indicate very high similarity, e.g., the group (\tilde{y} , i , η).

4. First Decipherment Approach

Building on the self-similarity of Roman letters, our first theory was that the Roman letters carry all the information in the cipher, and that all other symbols are NULLs (meaningless tokens added after encipherment to confuse cryptanalysis). If we remove all other symbols, the remaining Roman letters indeed follow a typical natural language distribution, with the most popular letter occurring 12% of the time, and the least popular letters occurring rarely.

The revealed sequence of Roman letters is itself nonsensical, so we posited a simple substitution cipher. We carried out automatic computer attacks against the revealed Roman-letter sequence, first assuming German source, then English, then Latin, then forty other candidate European and non-European languages. The attack method is given in [Knight et al, 2006]. That method automatically combines plaintext-language identification with decipherment. Unfortunately, this failed, as no

language identified itself as a more likely plaintext candidate than the others.

We then gave up our theory regarding NULLs and posited a homophonic cipher, with each plaintext letter being encipherable by any of several distinct cipher letters. While a well-executed homophonic cipher will employ a flat letter frequency distribution, to confound analysis, we guessed that the Copiale cipher is not optimized in this regard.

We confirmed that our computer attack does in fact work on a synthetic homophonic cipher, i.e., it correctly identifies the plaintext language, and yields a reasonable, if imperfect, decipherment. We then loosed the same attack on the Copiale cipher. Unfortunately, all resulting decipherments were nonsense, though there was a very slight numerical preference for German as a candidate plaintext language.

5. Second Decipherment Approach

We next decided to focus on German as the most likely plaintext language, for three reasons:

- the book is located in Germany
- the computer homophonic attack gave a very slight preference to German
- the book ends with the inscription “Philipp 1866”, using the German double-p spelling.

Pursuing the homophonic theory, our thought was that all five circumflexed letters (\hat{a} , \hat{e} , \hat{i} , \hat{o} , \hat{u}), behaving similarly, might represent the same German letter. But which German letter? Since the circumflexed letters are preceded by \mathfrak{z} and π , the circumflexed letters would correspond to the German letter that often follows *whatever z and pi stand for*. But what do they, in turn, stand for?

From German text, we built a digraph frequency table, whose the most striking characteristic is that C is almost always followed by H. The German CH pair is similar to the English QU pair, but C is fairly frequent in German. A similar digraph table for the cipher letters shows that \mathfrak{c} is almost always followed by \mathfrak{h} . So we posited our first two substitutions: $\mathfrak{c}=\text{C}$ and $\mathfrak{h}=\text{H}$. We then looked for what typically precedes and follows CH in German, and what typically precedes and follows $\mathfrak{c}\mathfrak{h}$ in the cipher.

For example, $\mathfrak{c}\mathfrak{h}\mathfrak{a}$ is the most frequent cipher trigraph, while CHT is a common German trigraph. We thus hypothesized the further substitution $\mathfrak{a}=\text{T}$, and this led to a cascade of others. We retracted any hypothesis that resulted in poor German digraphs and trigraphs, and in this way, we could make steady progress (Figure 5).

The cluster map in Figure 4 was of great help. For example, once we established a substitution like $\mathfrak{y}=\text{I}$, we could immediately add $\mathfrak{u}=\text{I}$ and $\mathfrak{v}=\text{I}$, because the three cipher letters behave so similarly. In this way, we mapped all circumflexed letters (\hat{a} , \hat{e} , \hat{i} , \hat{o} , \hat{u}) to plaintext E. These leaps were frequently correct, and we soon had substitutions for over 50 cipher letters.

Despite progress, some very frequent German trigraphs like SCH were still drastically under-represented in our decipherment. Also, many cipher letters (including all unaccented Roman letters) still lacked substitution values. A fragment of the decipherment thus far looked like this (where “?” stands for an as-yet-unmapped cipher letter):

?GEHEIMER?UNTERLIST?VOR?DIE?GESELLE
?ERDER?TITUL
?CEREMONIE?DER?AUFNAHME

On the last line, we recognized the two words CEREMONIE and DER separated by a cipher letter. It became clear that *the unaccented Roman letters serve as spaces in the cipher*. Note that this is the opposite of our first decipherment approach (Section 4). The non-Roman letters are not NULLs -- they carry virtually all the information. This also explains why paragraphs start with capitalized Roman letters, which look nice, but are meaningless.

We next put our hypothesized decipherment into an automatic German-to-English translator (www.freetranslation.com), where we observed that many plaintext words were still untranslatable. For example, ABSCHNITL was not recognized as a translatable German word. The final cipher letter for this word is colon (:), which we had mapped previously to L. By replacing the final L in ABSCHNITL with various letters of the alphabet (A-Z), we hit on the recognized word

αΠΔβ+ιΑγ.ζββϋθζεζγςήñππυ
 hyp: dos mit der andern hand
 corr: doch mit der andern hand

ρζιηλ:ίγ|εϋ̄ñî̄ξ̄η+:αβλὸρñcδρ+ιλφπυg.ζή=+υc
 hyp: dritlens einer n mlt tobach mit de daume
 corr: drittens einer ???? tobach mit dem daumen

κ̄ππ(ζυt+ιλ:δ̄c̄|υc̄|̄γ̄ḡχ̄ζsζρ̄h̄c̄ιρ̄δ̄ῡδ̄ή̄ñ̄z̄b̄
 hyp: und de mittelede finger der linche hand
 corr: und dem mittelsten finger der linchen hand

ρû|η̄ή̄r̄c̄s+ιλg.π̄ξ̄|δ̄c̄η̄ῡδ̄ῡl̄ή̄h̄ρ̄π̄m̄z̄β̄γ̄r̄ḡ
 hyp: beruhe mit der linche hand dein
 corr: berühre mit der linchen hand dein

This allowed us to virtually complete our table of substitutions (Figure 6). Three cipher letters remained ambiguous:

- ϩ could represent either SS or S
- Ϙ could represent either H or K
- υ could represent either EN or EM

However, these three symbols are ambiguous only with respect to deciphering into modern German, not into old German, which used different spelling conventions.

The only remaining undeciphered symbols were the large ones: ϰ, Ϡ, Δ, ϫ, ϙ, ϰ, ϰ, ϰ, and Π. These appear to be logograms, standing for the names of (doubly secret) people and organizations, as we see in this section: “the ϰ asks him whether he desires to be ϙ”.

6. Contents

The book describes the initiation of “DER CANDIDAT” into a secret society, some functions of which are encoded with logograms. Appendix A contains our decipherment of the beginning of the manuscript.

7. Conclusion

We described the Copiale cipher and its decipherment. It remains to transcribe the rest of the manuscript and to do a careful translation. The document may contain further encoded information, and given the amount of work it represents, we believe there may be other documents using the same or similar schemes.

Plaintext (German)	Ciphertext
A	î ñ ã ϣ*
Ä	ϣ*
B	β
C	γ
D	π ζ
E	â ê î ô û ϩ ϫ
F	ϣ
G	δ χ
H	η Ϙ*
I	ÿ η ι
J	τ
K	Ϙ*
L	ç
M	+
N	π ρ υ g
O	Δ ò
Ö	ϣ
P	δ
R	ř ç i
S	ϩ*
T	^
U	= ϫ
Ü	ϣ
V	δ
W	ñ
X	ƒ
Y	∞
Z	š
SCH	†
SS	ϩ*
ST	†
CH	†
repeat previous consonant	:
EN / EM	υ
space	a b c d e f g h i k l m n o p q r s t u v w x y z

Figure 6. Letter substitutions resulting from decipherment. Asterisks (*) indicate ambiguous cipher letters that appear twice in the chart. This table does not include the large characters: ϰ, Ϡ, Δ, ϫ, ϙ, ϰ, ϰ, and Π.

8. Acknowledgments

Thanks to Jonathan Graehl for plotting the cluster map of cipher letters. This work was supported in part by NSF grant 0904684.

9. References

K. Knight, A. Nair, N. Rathod, and K. Yamada, "Unsupervised Analysis for Decipherment Problems", Proc. ACL-COLING, 2006.

Appendix A

Ciphertext:

v x 3 | i l a s k p k j w n
 π ο ι ρ η Δ ρ ε ζ ε α = ρ λ u b ◊ u r ⊙ ε
 ϑ η η ζ ι + ϑ ι η ρ λ η i η ε ρ .
 c u i j ε z t p p t g η λ : κ
 ϑ x u η ε η + x r p k m λ i z : γ ρ λ ϑ ο i q z i x a i ε c : u ϑ .
 f u r i ϑ i κ λ i λ = c ε .
 m ρ η r a + Δ g y u x z ϑ z m n k Γ η h η + i f .
 κ m u r : p z i ϑ f j y o η ε i η ζ i λ n π a z b Δ g z = j ρ l z u ρ c λ a r g κ λ η
 η r η η λ x i η ρ i Δ r ϑ η λ a = g z w π y e c Δ r ϑ Δ + b z η r i x y i r z z u f λ x
 π x i ρ ϑ r = Γ i u l a s k m ϑ m i η g a i κ η = λ η i l x ϑ ρ f : r i λ b i u m y i z v
 z a i x o p m π i z h λ c ϑ o g g z u + r k η η x ε z g h λ η η i η η λ i z t a k i r ϑ y
 m a + h h r z ϑ z η b s η + : x r k r ρ ϑ η η Δ c u λ g = u z κ p x o o n z ρ z f η n r π
 x ε y u g = r π g ϑ Δ z n z η i κ π η i l x y r i g π u λ b g λ i t b ϑ u f η η z ϑ λ e
 z η ϑ i u r c Γ z ϑ ϑ λ b η η m ϑ :
 n i r i c i ρ e ρ h ϑ r ϑ r x ϑ u η z a k ◊ g s = ϑ m ε j z u l i
 g s m u o o λ ϑ η | o π ε g u ϑ ϑ x r Δ j z g x r x u g π u z κ ⊙ η i η ρ b = u
 λ ε r m ϑ z f : u a = r z l Δ η r η ϑ m η a z ϑ i j d i r i ϑ ϑ y λ u z i ε g c u η
 r s ε η λ
 c f = i f a η u b m Δ c : x ϑ .
 h z i η λ : ϑ g | e z η a x Δ n π a z ⊙ p s k r ϑ i z t m a y ϑ u l i = n π
 z p s = n h k f r z p i n ϑ x i m g π c y o η f ϑ b i ρ κ a g η x y t i ϑ c s = b + n ρ u ϑ
 x ϑ i r : u λ v i o o n .
 i π ϑ z f o h g z y p r a l i m λ m Δ r λ ρ o e q n r .

Plaintext, as deciphered:

gesetz buchs
 der hocherleuchte ◊ e ⊙
 geheimer theil.
 erster abschnitt
 geheimer unterricht vor die gesellen.
 erster titul.
 ceremonien der aufnahme.
 wenn die sicherheit der Δ durch den ältern
 thürhuter besorget und die Δ vom dirigirenden λ
 durch aufsetzung seines huths geöffnet ist wird der
 candidat von dem jüngern thürhüter aus einem andern
 zimmer abgehohlet und bey der hand ein und vor des
 dirigirenden λ tisch geführet dieser frägt ihn:
 erstlich ob er begehre ◊ zu werden
 zweytens denen verordnungen der ⊙ sich
 unterwerffen und ohne widerspenstigkeit die lehrzeit
 ausstehen wolle.
 drittens die Δ der ⊙ gu verschweigen und dazu
 auf das verbindlichste sich anheischig zu machen
 gesinnet sey.
 der candidat antwortet ja.

Initial translation:

First lawbook
 of the ◊ e ⊙
 Secret part.
 First section
 Secret teachings for apprentices.
 First title.
 Initiation rite.
 If the safety of the Δ is guaranteed, and the Δ is
 opened by the chief λ, by putting on his hat, the
 candidate is fetched from another room by the
 younger doorman and by the hand is led in and to the
 table of the chief λ, who asks him:
 First, if he desires to become ◊.
 Secondly, if he submits to the rules of the ⊙ and
 without rebelliousness suffer through the time of
 apprenticeship.
 Thirdly, be silent about the Δ of the ⊙ and
 furthermore be willing to offer himself to volunteer
 in the most committed way.
 The candidate answers yes.