



OCR of Legacy Documents

as a Building Block in Industrial Disaster Prevention

Isemann, D., Niekler, A., Preßler, B., Viereck F. and Heyer G.



Background

Project: “Knowledge Management for Legacy Documents in Science, Administration and Industry”

The Asse Pit



Asse History



- 1909-1964: Salt mine
- 1964-1978: “Exploratory final storage” (Versuchsendlager) for radioactive waste
- 1. 1. 1979: No new radioactive material accepted
- 1979—now: Managing and maintaining the site and its contents
- 1988: Ingress of water into the pit
- 2009: Change of management, political pressures

Asse Legacy Files



- In the order of 5000 file folders documenting / including (no claim to completeness):
 - Geological and mining surveys
 - Radioactive waste management and reactive transport studies
 - Scientific and administrative correspondence
 - Telephone protocols and meeting minutes
 - Delivery notes / receipts for radioactive substances
 - Reports / assessments of extraordinary events (e.g. earthquakes, operational incidents, etc.)
 - Guidelines and written exhortations making reference to legal and regulatory frameworks governing nuclear waste disposal
- Huge source of information of political, historical and practical interest (current plan is to empty out the nuclear waste repository)

Example Documents

**GESELLSCHAFT FÜR KERNFORSCHUNG M.B.H.
KARLSRUHE**

Gesellschaft für Kernforschung m. b. H. - 75 Karlsruhe - Postfach 947

An das
Bundesministerium für
Wissenschaftliche Forschung
532 Bad Godesberg
Luisenstraße 46

Der Bundesminister für
Wissenschaftliche Forschung
am 24. FEB. 1965
75 KARLSRUHE
ANL. 3
AZ. 1875-3-79/6
WELLESSTRASSE 3

IHRER ZEICHEN: / IHRER NACHRICHT VOM: / UNSERE ZEICHEN: / TAG: / MAUSNUM: /

Dr. H./RZ 2010 19.2.1965

BESPRECHUNG ÜBER GRUNDWASSERGUTACHTEN FÜR SCHACHTANLAGE ASSE

Sehr geehrte Herren!

Am 12.2.1965 fand gemäß einer telefonischen Aufforderung Ihres Hauses vom 9.2.d.J. eine Besprechung mit Herrn Prof. Semmler, Westfälische Berggewerkschaftskasse, Abteilung Wasserwirtschaft und Hydrogeologie, über "Grundwassergutachten für Schachtanlage Asse" statt. Eine Notiz darüber liegt bei.

Herr Prof. Semmler ist bereit, ein Gutachten mit dem Titel "Geführungsmöglichkeit der Trinkwasserversorgung in der Umgebung des Kalibergwerkes Asse II durch Einlagerung von radioaktiven Abfällen in aufgelassenen Grubenbauen" anzufertigen. Die weiteren grundlegenden Untersuchungen können nur von der Westfälischen Berggewerkschaftskasse, Abteilung Wasserwirtschaft, in ihrer Gesamtheit durchgeführt werden. Für einen Auftrag gilt die beiliegende Gebührenordnung.

Wir bitten um Ihr Einverständnis zu Aufträgen an

- Herrn Prof. Semmler (persönlich) über die Erstellung eines Vorgutachtens mit nachfolgendem Titel "Geführungsmöglichkeiten der Trinkwasserversorgung in der Umgebung des Kalibergwerkes Asse II durch Einlagerung von radioaktiven Abfällen in den aufgelassenen Grubenbauen"

- 2 -

24. FEB. 1965

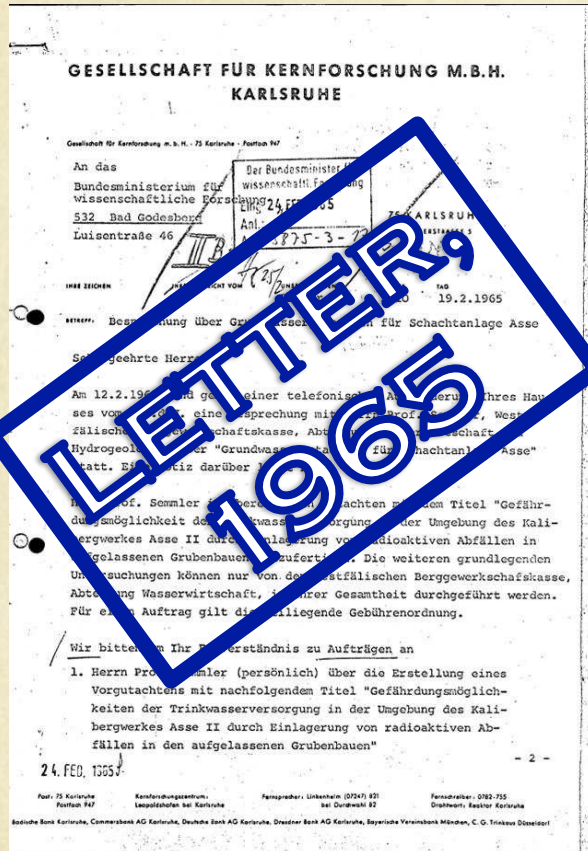
Post: 75 Karlsruhe Postfach 947
Kernforschungszentrum, Langenfelder Str. Karlsruhe
Fernsprecher: Linienamt (076) 821 840, Durchwahl 82
Fernschreiber: 0762-755
Drahtfunk: Essener Karlsruhe
Bayerische Bank Karlsruhe, Commerzbank AG Karlsruhe, Deutsche Bank AG Karlsruhe, Dresdner Bank AG Karlsruhe, Bayerische Vereinsbank München, C. G. Tinkaus Düsseldorf

INHALTSVERZEICHNIS

	Seite
I. EINFÜHRUNG	1
II. ARBEITSMETHODIK	3
III. GEOMORPHOLOGISCHE ÜBERSICHT	4
IV. STRATIGRAPHIE	7
IV.1. Zechstein	7
IV.2. Buntsandstein	10
IV.3. Muschelkalk	24
IV.4. Keuper	38
IV.5. Lias und Dogger	47
IV.6. Kreide (Unter-Hauterivium)	50
IV.7. Tertiär (Unter-Oligozän)	52
IV.8. Quartär	60
V. TEKTONIK	70
V.1. "Bereich des verstärzten Deckgebirges"	70
V.2. Der tektonische Bau des Salzgebirges	72
V.2.1. Schacht Asse II	73
V.2.2. Schacht Asse I	79
V.3. Der tektonische Bau des Deckgebirges	81
V.3.1. Art und Entstehung der tektonischen Formen	81
V.3.1.1. Verwerfungen	81
V.3.1.2. Klüftung	86
V.3.1.3. Tangentiale Einengung?	88
V.3.2. Die räumliche Verteilung der tektonischen Formen und ihre Bedeutung für die Strukturbildung	90
V.3.2.1. Der nordwestliche Strukturschluß und der Salzstock von Gr. Denkte	90
V.3.2.2. Die Asse i.e.S.	97
V.3.2.2.1. Der NW-Abschnitt der SW-Flanke	97
V.3.2.2.2. Der SE-Abschnitt der SW-Flanke	101
V.3.2.2.3. Der Mittelabschnitt der SW-Flanke	103
V.3.2.3. Die südöstliche Verlängerung der Asse	108
V.4. Der zeitliche Ablauf der Strukturbildung	114
VI. ZUSAMMENFASSUNG	118
VII. LITERATURVERZEICHNIS	122
VIII. ANHANG	
geologische Profile	
Streichlinienkarte der SW-Flanke	
Karte der Quartärverbreitung und	
geologische Karte	-mächtigkeit



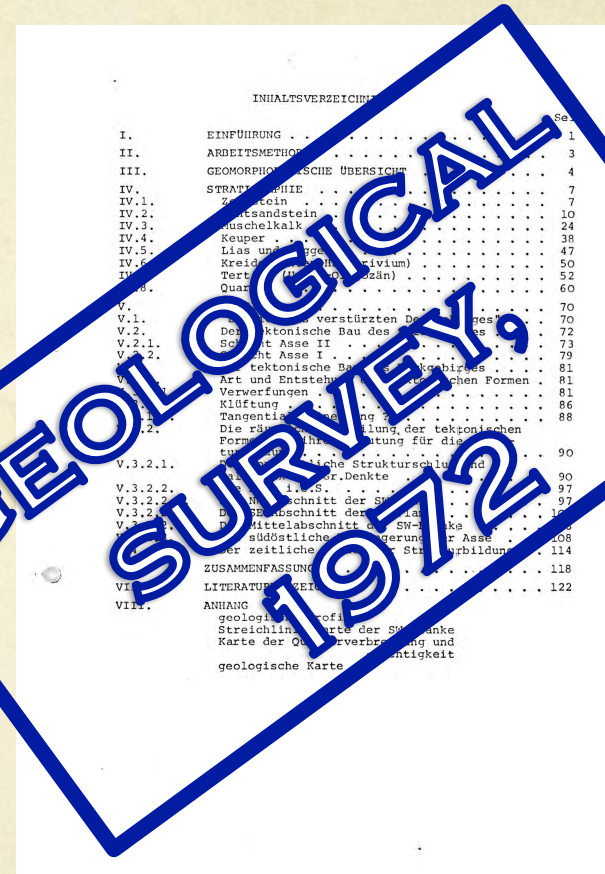
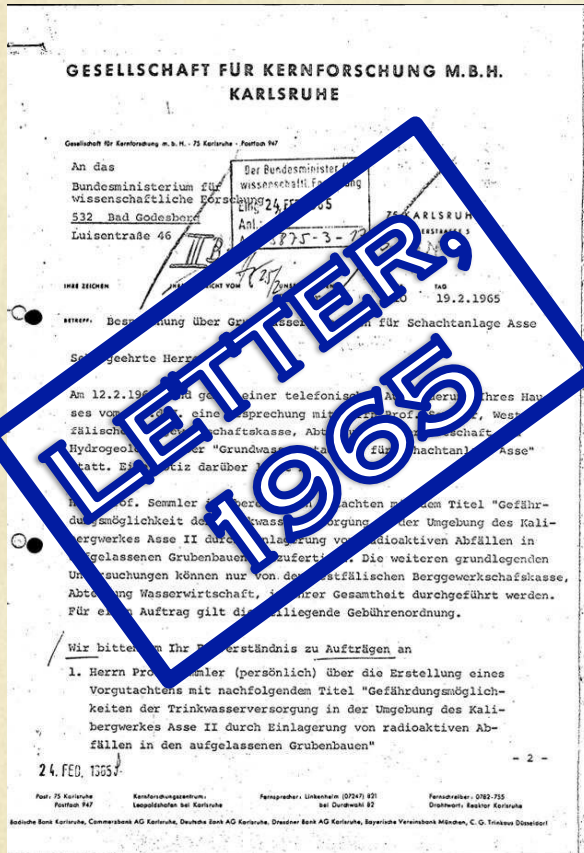
Example Documents



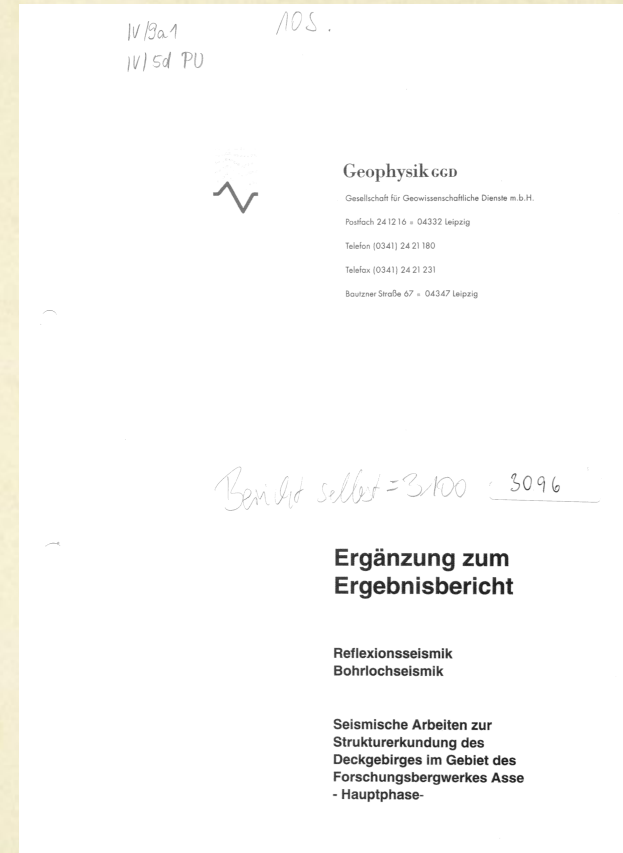
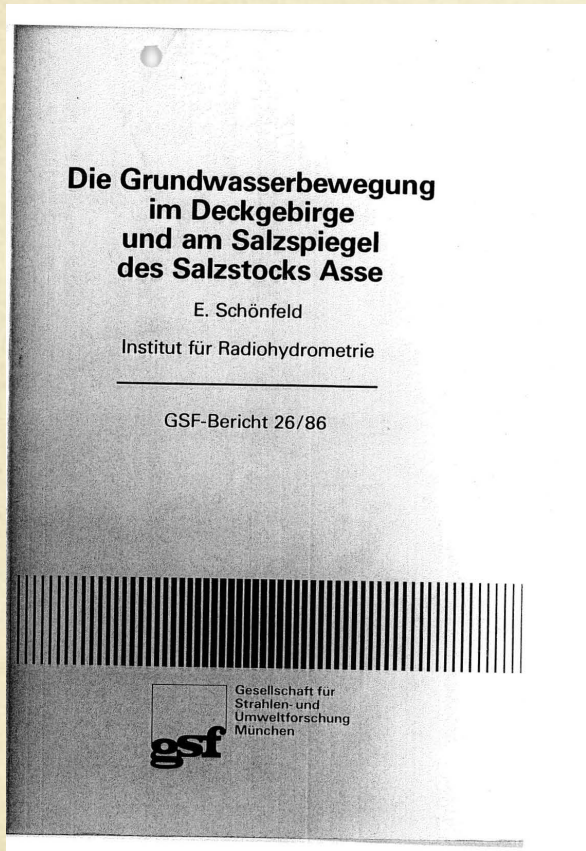
INHALTSVERZEICHNIS

	Seite
I.	EINFÜHRUNG 1
II.	ARBEITSMETHODIK 3
III.	GEOMORPHOLOGISCHE ÜBERSICHT 4
IV.	STRATIGRAPHIE 7
IV.1.	Zechstein 7
IV.2.	Buntsandstein 10
IV.3.	Muschelkalk 24
IV.4.	Keuper 38
IV.5.	Lias und Dogger 47
IV.6.	Kreide (Unter-Hauterivium) 50
IV.7.	Tertiär (Unter-Oligozän) 52
IV.8.	Quartär 60
V.	TEKTONIK 70
V.1.	"Bereich des verfestigten Deckgebirges" 70
V.2.	Der tektonische Bau des Salzgebirges 72
V.2.1.	Schacht Asse II 73
V.2.2.	Schacht Asse I 79
V.3.	Der tektonische Bau des Deckgebirges 81
V.3.1.	Art und Entstehung der tektonischen Formen 81
V.3.1.1.	Verwerfungen 81
V.3.1.2.	Klüftung 86
V.3.1.3.	Tangentiale Einengung ? 88
V.3.2.	Die räumliche Verteilung der tektonischen Formen und ihre Bedeutung für die Strukturbildung 90
V.3.2.1.	Der nordwestliche Strukturschluß und der Salzstock von Gr.Denkte 90
V.3.2.2.	Die Asse i.e.S. 97
V.3.2.2.1.	Der NW-Abschnitt der SW-Flanke 97
V.3.2.2.2.	Der SE-Abschnitt der SW-Flanke 101
V.3.2.2.3.	Der Mittelabschnitt der SW-Flanke 103
V.3.2.3.	Die südöstliche Verlängerung der Asse 108
V.4.	Der zeitliche Ablauf der Strukturbildung 114
VI.	ZUSAMMENFASSUNG 118
VII.	LITERATURVERZEICHNIS 122
VIII.	ANHANG
	geologische Profile
	Streichlinienkarte der SW-Flanke
	Karte der Quartärverbreitung
	geologische Karte -mächtigkeit

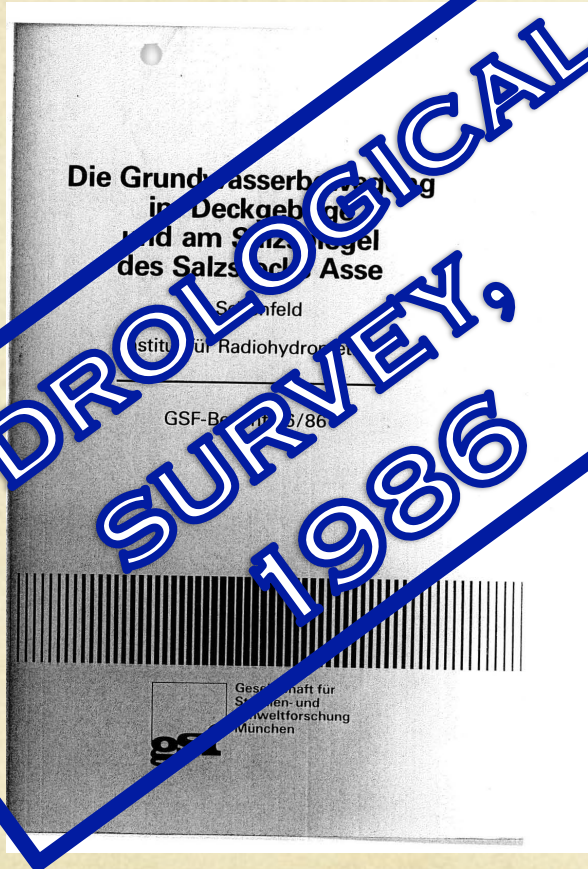
Example Documents



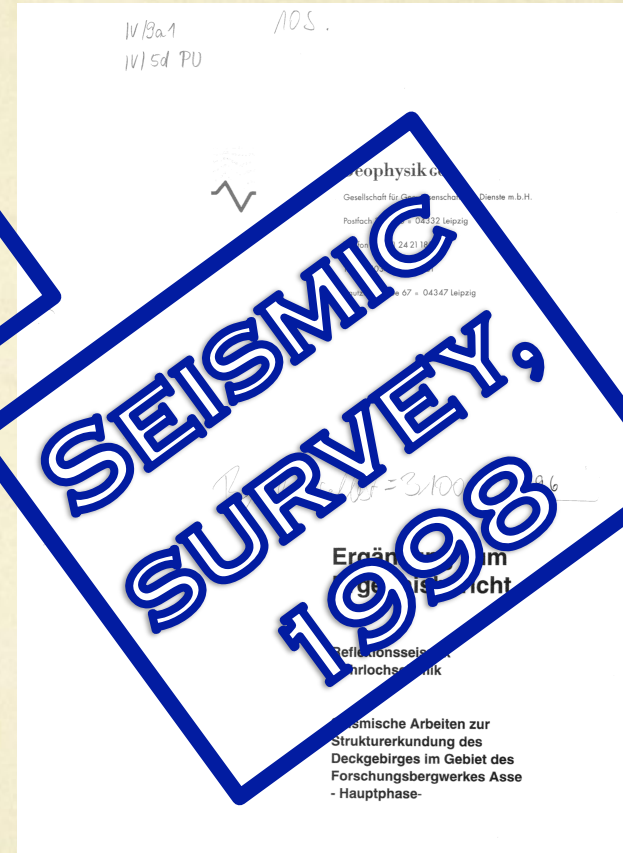
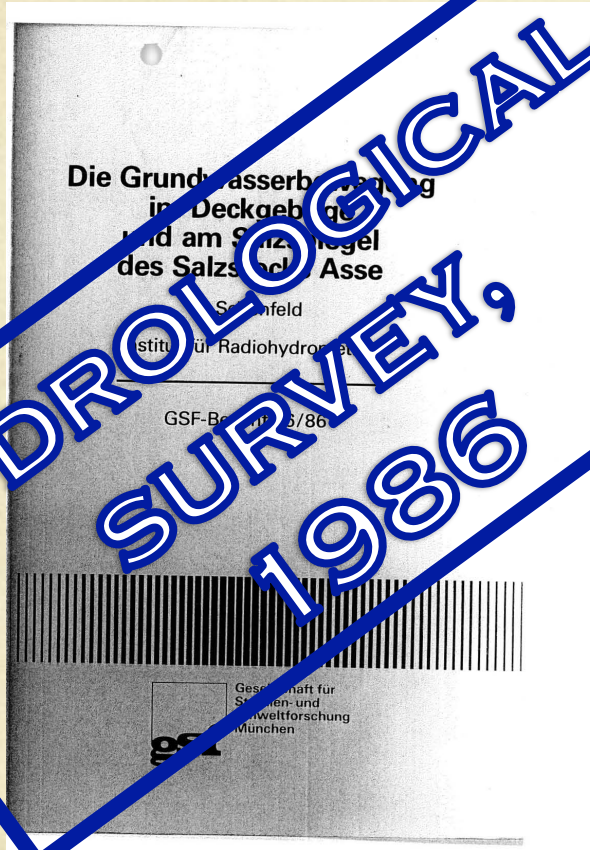
Example Documents



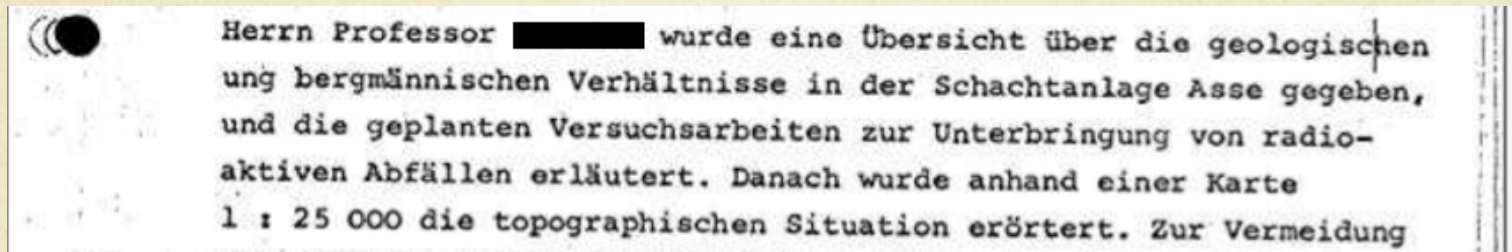
Example Documents



Example Documents



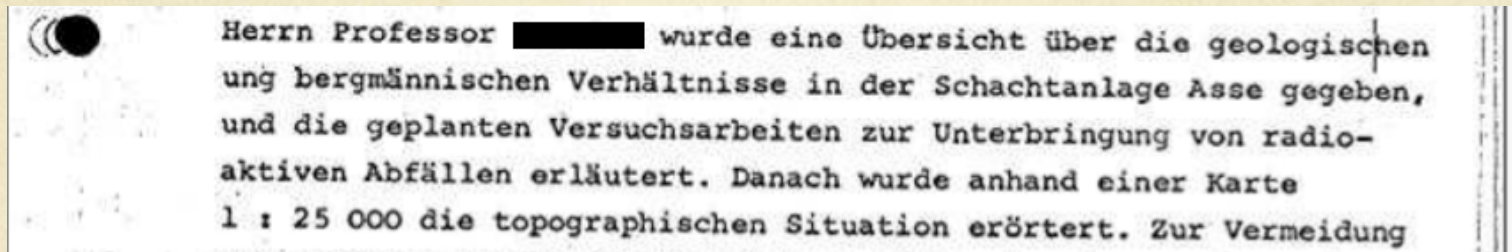
Challenge for OCR



Google Drive OCR

. | Herrn Professor XXXXX wurde eine: Über:=icht über? du: genlug:..scfuan um; hergmännischen Uerhältnisse in der: Schachtanlage ,um: und die qaplanten versuchsarbaitan zu: Unterbringung man radia MEÄILEn erläutert. Danach wurde anhand einer Harta 1 : 25 üüü dia top-ugraphischen Eituatinn arä:terr:.. Zur vermeiflung

Challenge for OCR



Google Drive OCR (correctly identified tokens)

. | Herrn Professor XXXXX wurde eine: Über:=icht über? du: genug:..scfuan
um; hergmännischen Uerhältnisse in der: Schachanlage ,um:
und die qaplanten versuchsarbaitan zu: Unterbringung man radia
MEÄILEn erläutert. Danach wurde anhand einer Harta
1:25 üüü dia top-ugraphischen Eitutinn arä:terr:.. Zur vermeiflung

Challenge for OCR

Herrn Professor [REDACTED] wurde eine Übersicht über die geologischen und bergmännischen Verhältnisse in der SchachanlageASSE gegeben, und die geplanten Versuchsarbeiten zur Unterbringung von radioaktiven Abfällen erläutert. Danach wurde anhand einer Karte 1 : 25 000 die topographische Situation erörtert. Zur Vermeidung

Google Drive OCR (correctly identified tokens)

Herrn Professor wurde
in
und die
erläutert. Danach wurde anhand einer
1 : 25 Zur

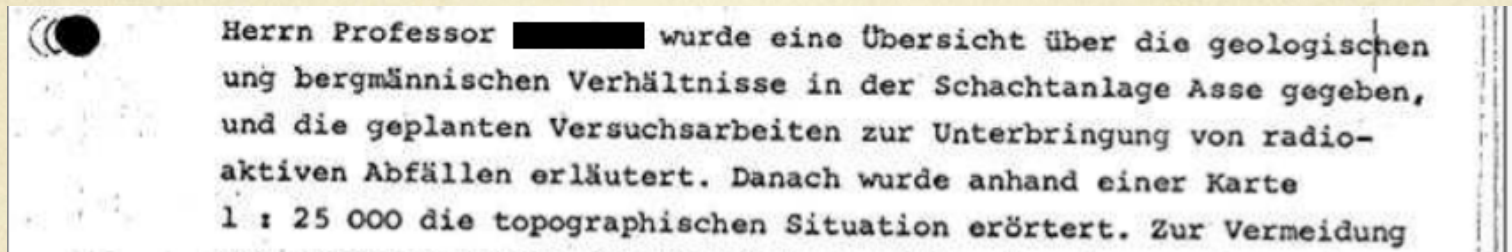
Challenge for OCR

Herrn Professor [REDACTED] wurde eine Übersicht über die geologischen und bergmännischen Verhältnisse in der SchachtanlageASSE gegeben, und die geplanten Versuchsarbeiten zur Unterbringung von radioaktiven Abfällen erläutert. Danach wurde anhand einer Karte 1 : 25 000 die topographische Situation erörtert. Zur Vermeidung

Google Drive OCR (correctly identified tokens)

Herrn Professor [REDACTED] wurde
in
und die
erläutert. Danach wurde anhand einer
1 : 25 [REDACTED] Zur

Challenge for OCR

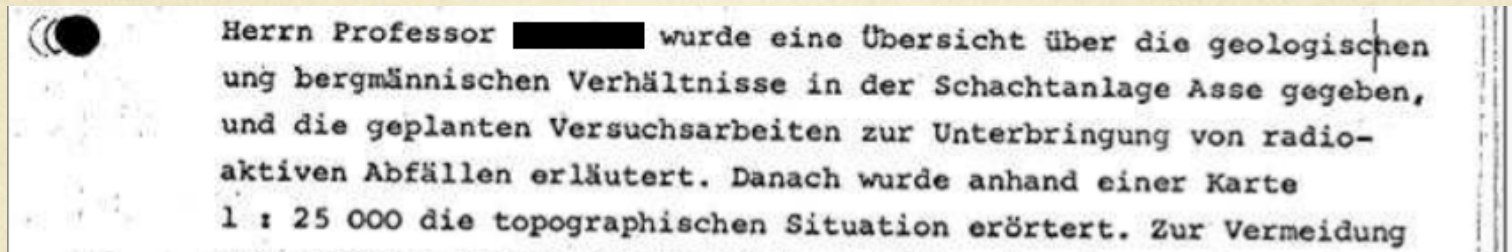


Online OCR service

Herrn Professor XXXXXX wurde ein* über die geologischen
14.11/4.9111111

une bergmännischen Verhältnisse in der Schachanlage Asse gegeben
und die geplanten Versuchsarbeiten zur Unterbringung von radio-
aktiven Abfällen erläutert. Danach wurde anhand einer Karte
1 : 25 000 die topographische Situation erörtert. Zur Vermeidung

Challenge for OCR

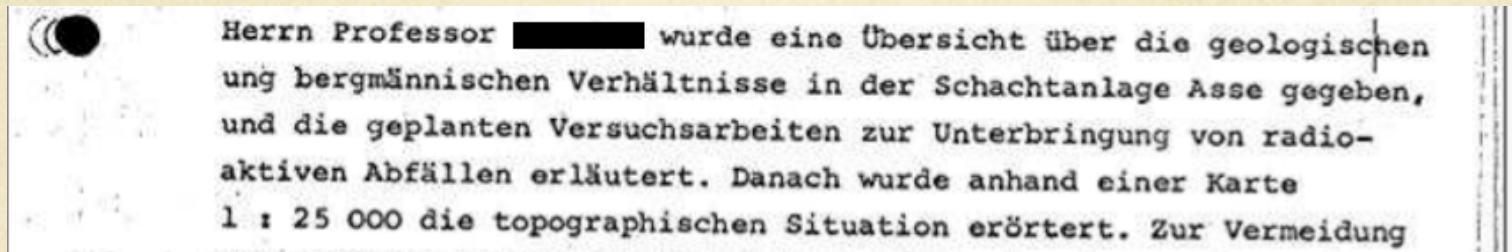


Online OCR service (correctly identified tokens)

Herrn Professor wurde über geologischen

 bergmännischen Verhältnisse in der Schachanlage Asse
und die geplanten zur von radio-
aktiven Abfällen erläutert. Danach wurde anhand einer Karte
1 : 25 000 die topographische Situation erörtert. Zur Vermeidung

Challenge for OCR



Online OCR service (correctly identified tokens)

Better, but
keywords
are missing

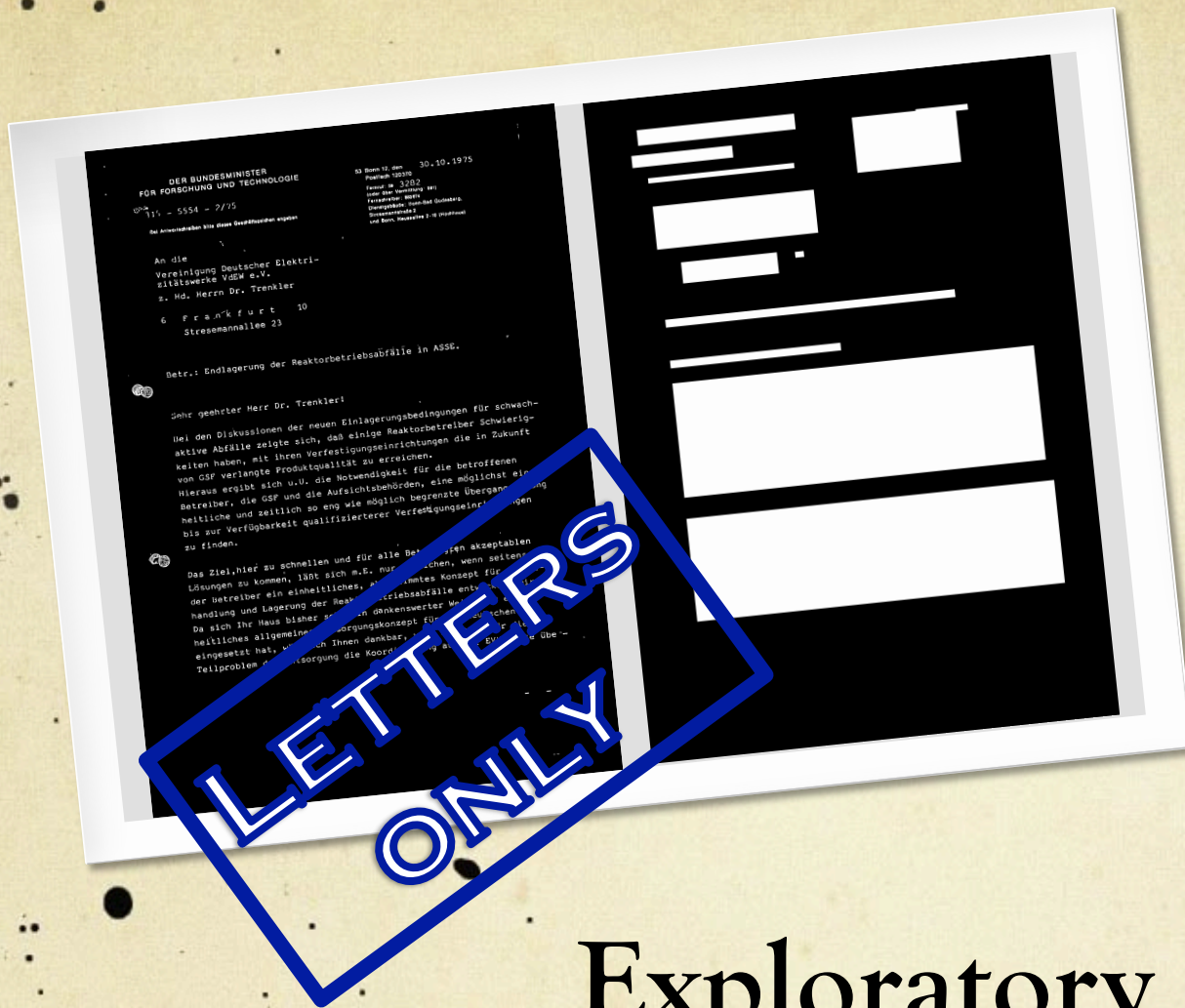
Herrn Professor wurde über geologischen

bergmännischen Verhältnisse in der SchachanlageASSE
und die geplanten zu von radio-
aktiven Abfällen erläutert. Danach wurde anhand einer Karte
1 : 25 000 die topographischen Situation erörtert. Zur Vermeidung



Vision for the Project

- Archiving: Searchable access to legacy paper documents
- Repurposing: Extracting metadata, meaningful connections, ontology learning
- Integrated processing chain from OCR to semantics
- First exploratory study: Segmentation of letterheads

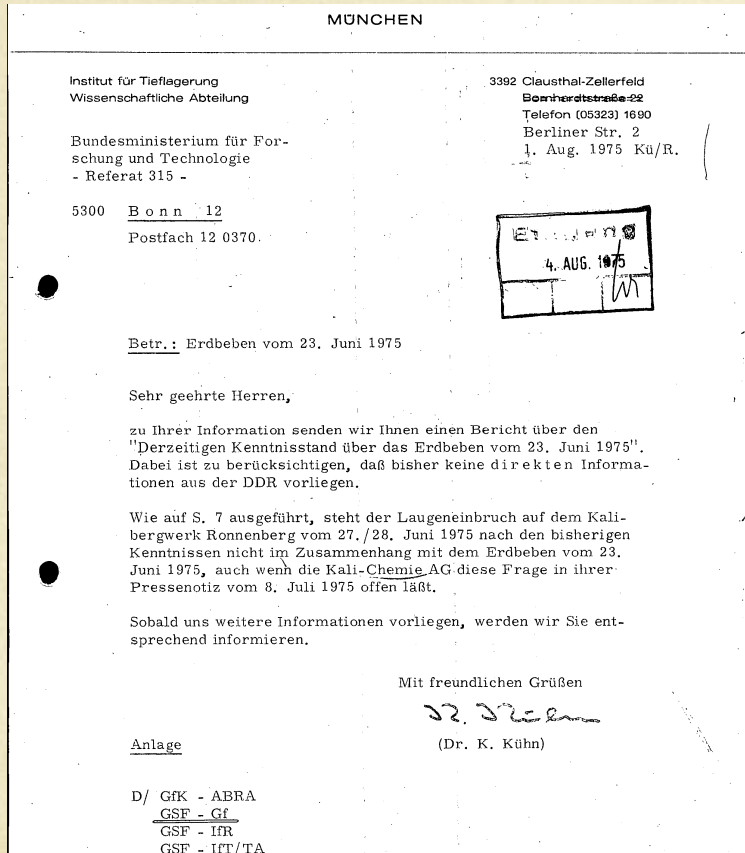


Exploratory Study

Letterhead segmentation for legacy correspondence

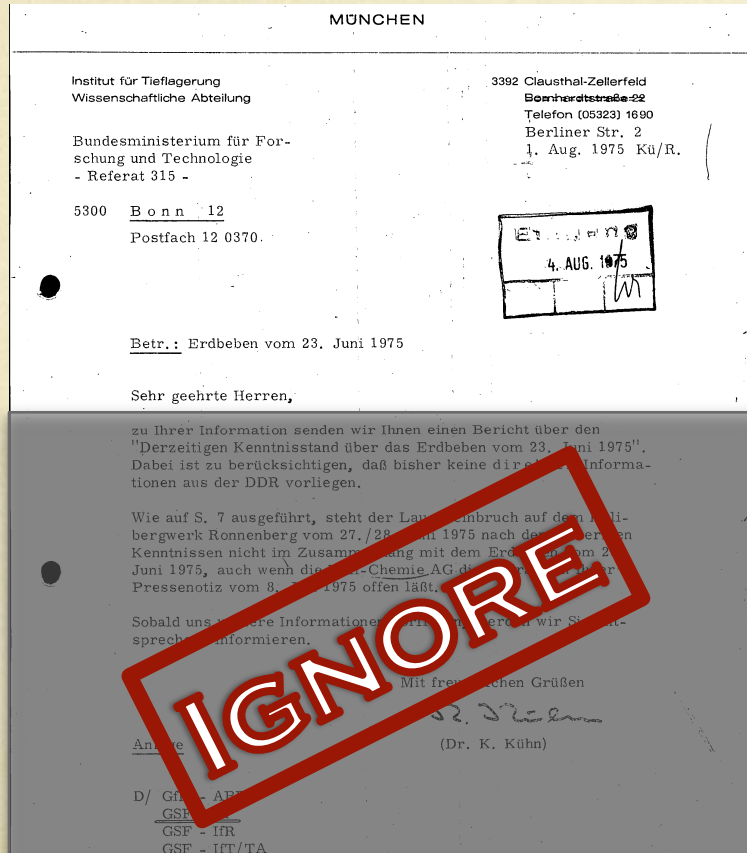
20

Data



- Want to identify letterhead elements that form a unit
- 24 letters from an initial batch of scanned documents (approx. 80 documents)
- OCR optimised by external project partner
- Text body stripped away, just letterheads

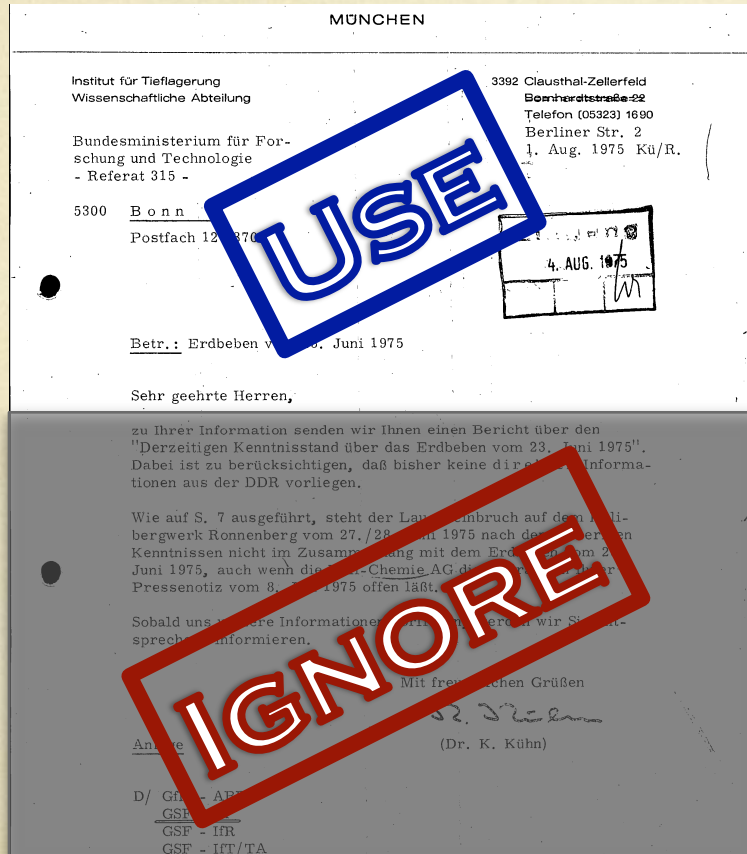
Data



- Want to identify letterhead elements that form a unit
- 24 letters from an initial batch of scanned documents (approx. 80 documents)
- OCR optimised by external project partner
- Text body stripped away, just letterheads

22

Data



- Want to identify letterhead elements that form a unit
- 24 letters from an initial batch of scanned documents (approx. 80 documents)
- OCR optimised by external project partner
- Text body stripped away, just letterheads

23

Method

```
<span class='ocr line'  
      title='bbox  
      261 2648 1064 2719'  
      style=  
      'position:absolute;  
left:93.2142857143px;  
top:313.928571429px'>  
      Forschung und</  
      span><br />
```

Fig.: Snapshot from hOCR output

- Unsupervised clustering approach
- Use only positional information (for now)
- Agglomerative, bottom-up clustering
 - Complete-linkage
 - Euclidian distance

Method

```
<span class='ocr line'  
      title='bbox  
      261 2648 1064 2719'  
      style=  
      'position:absolute;  
left:93.2142857143px;  
top:313.928571429px'>  
      Forschung und</  
      span><br />
```

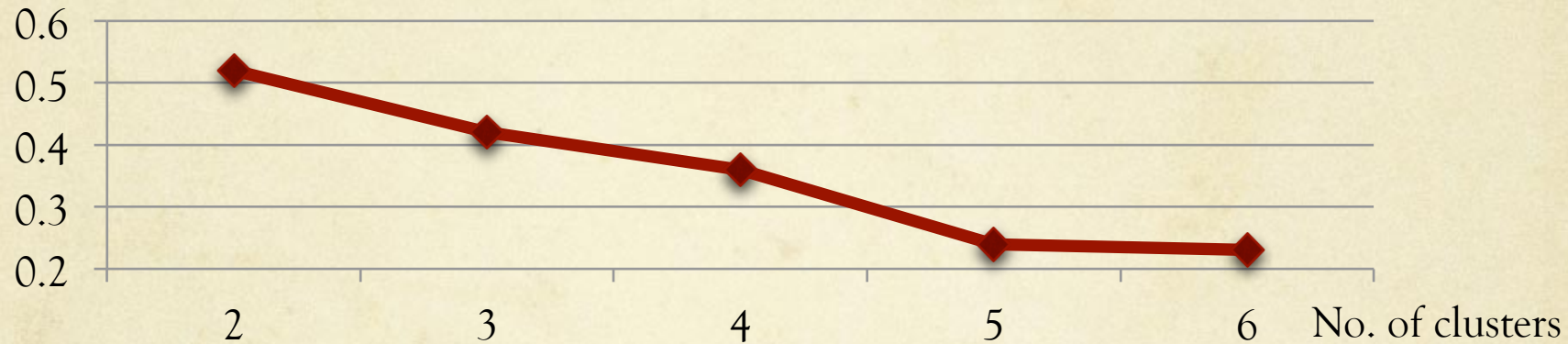
Fig.: Snapshot from hOCR output

- Unsupervised clustering approach
- Use only positional information (for now)
- Agglomerative, bottom-up clustering
 - Complete-linkage
 - Euclidian distance

Method

Avg. within
cluster
variance

Doc. #3



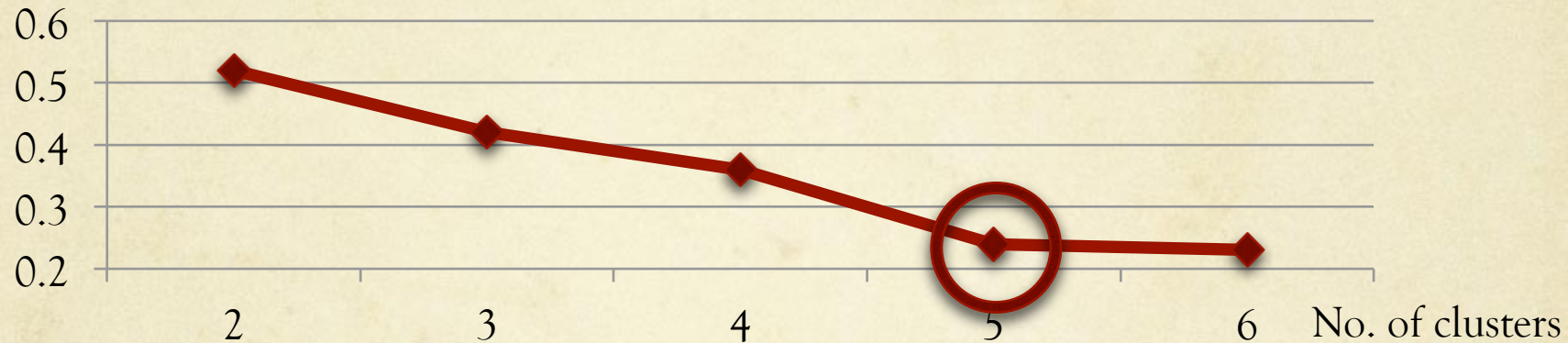
- “Flatten” clustering by identifying significant increase in average within cluster variance (“spot the elbow joint”)
- This is likely to indicate that two clusters far away from each other have erroneously been merged

26

Method

Avg. within
cluster
variance

Doc. #3

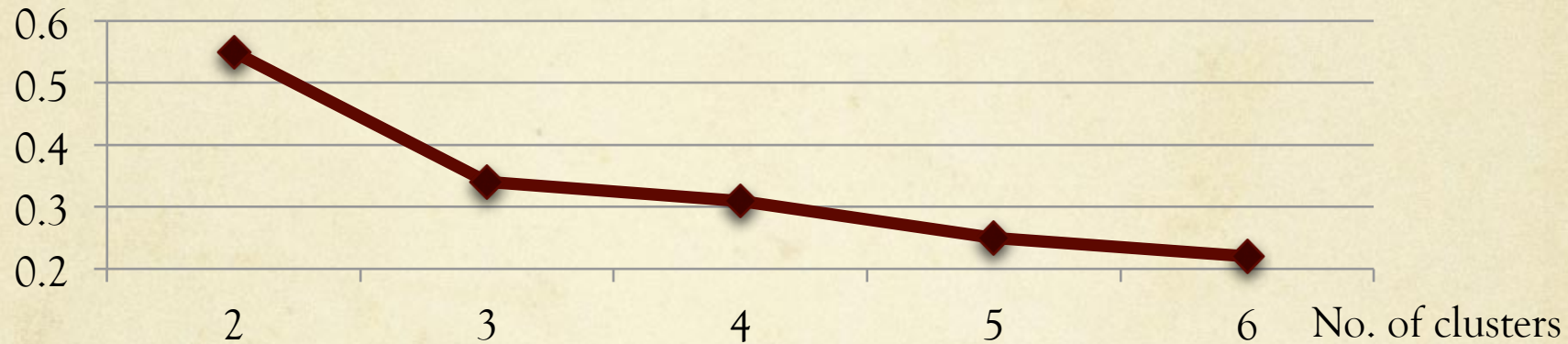


- “Flatten” clustering by identifying significant increase in average within cluster variance (“spot the elbow joint”)
- This is likely to indicate that two clusters far away from each other have erroneously been merged

Method

Avg. within
cluster
variance

Doc. #54



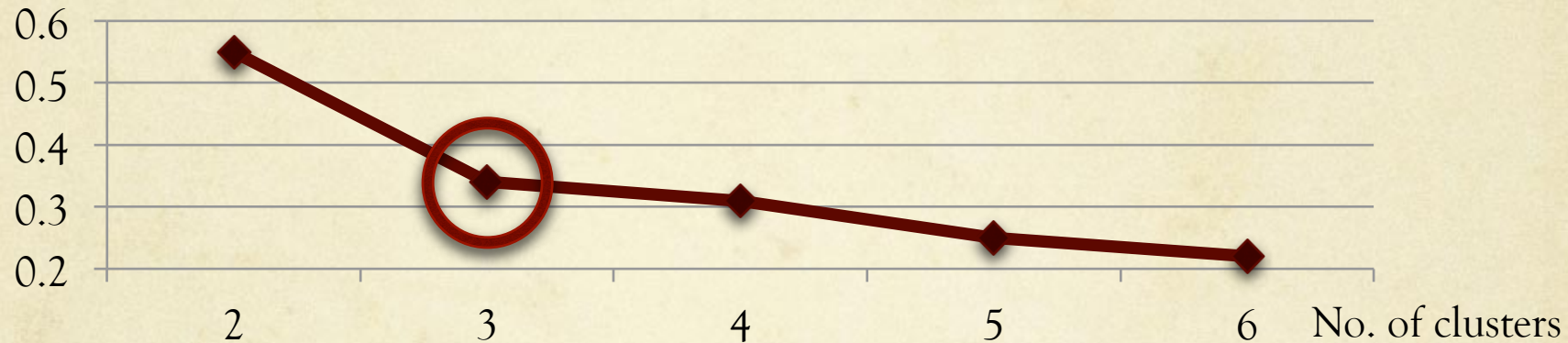
- “Flatten” clustering by identifying significant increase in average within cluster variance (“spot the elbow joint”)
- This is likely to indicate that two clusters far away from each other have erroneously been merged

28

Method

Avg. within
cluster
variance

Doc. #54



- “Flatten” clustering by identifying significant increase in average within cluster variance (“spot the elbow joint”)
- This is likely to indicate that two clusters far away from each other have erroneously been merged

29

Example Segmentation

BUNDESMINISTER
FOR FORSCHUNG UND TECHNOLOGIE

31'i - 5554 - 7/15

Tel Anwortehrelben Blick Degen Renditezielen angeben

An die

Vereinigung Deutscher
Elektrizitätswerke VdEW e.V.
z. Hd. Herrn Dr. XXXXX XXX

F r a n k f u r t
Stresemannallee 23

P-- A--. 7A)

Petr. Endlagerung der Reaktorbetriebsabfälle in ASSE.

ehr geehrter Herr Dr. XXXXXXXXX

A. Je-t Wu New'

53 Bonn 12, den 30.10.1975

Konery sw 3282
(oder aber Vermittlung 591)

Dienstgebäude: Bonn-Bad Godesberg.
und Bonn, Heuballen 2- 10 (Wohnung)

No. of
clusters

Within
cluster
variance

Example Segmentation

OER BUNDESMINISTER
 FOR FORSCHUNG UND TECHNOLOGIE
 31'i - 5554 - 7/15
 Tel Anwortehrelben Blick Degen Renditezielen angeben

53 Bonn 12, den 30.10.1975
 Koneu sw 3282
 (oder aber Vermittlung 591)
 Dienstgebv\$ude: Bonn-Bad Godesberg.
 und Bonn, Heuballen 2- 10 (Wohnung)

An die
 Vereinigung Deutscher
 Elektriziv\$ts werke VdEW e.V.
 z. Hd. Herrn Dr. XXXXX XXX
 F r a n k f u r t
 Stresemannallee 23

Petr. Endlagerung der Reaktorbetriebsabfv\$11e in ASSE.

chr geehrter Herr Dr. XXXXXXXXX

P-- A--. 7A)

A. Je-t Wu New'

No. of clusters	Within cluster variance
6	0.23

Example Segmentation

DER BUNDESMINISTER
 FOR FORSCHUNG UND TECHNOLOGIE
 31'i - 5554 - 7/15
 Tel Antwortehreiben Blick Degen Renditezielen angeben

53 Bonn 12, den 30.10.1975
 Koneu sw 3282
 (oder aber Vermittlung 591)
 Dienstgebv/sude: Bonn-Bad Godesberg.
 und Bonn, Heuballen 2- 10 (Wohnung)

An die
 Vereinigung Deutscher
 Elektrizit/swerke VdEW e.V.
 z. Hd. Herrn Dr. XXXXX XXX
 F r a n k f u r t
 Stresemannallee 23

Petr. Endlagerung der Reaktorbetriebsabfv/11e in ASSE.
 chr geehrter Herr Dr. XXXXXXXXX

P-- A--. 7A)
 A. Je-t Wu New'

No. of clusters	Within cluster variance
6	0.23
5	0.24

Example Segmentation

BUNDESMINISTER
 FÜR FORSCHUNG UND TECHNOLOGIE
 31'i - 5554 - 7/15
 Tel Anwortehreiben Blick Degen Renditezielen angeben
 An die
 Vereinigung Deutscher
 Elektrizitätswerke VdEW e.V.
 z. Hd. Herrn Dr. XXXXX XXX
 Frankfurt
 Stresemannallee 23

53 Bonn 12, den 30.10.1975
 Konev sw 3282
 (oder aber Vermittlung 591)
 Dienstgebäude: Bonn-Bad Godesberg.
 und Bonn, Heuballen 2- 10 (Wohnung)

Petr. Endlagerung der Reaktorbetriebsabfälle in ASSE.
 ehr geehrter Herr Dr. XXXXXXXXX

P-- A--. 7A)
 A. Je-t Wu New'

No. of clusters	Within cluster variance
6	0.23
5	0.24
4	0.36

Example Segmentation

DER BUNDESMINISTER
FÜR FORSCHUNG UND TECHNOLOGIE

31'i - 5554 - 7/15

Tel Anwortehrelben Blick Degen Renditezielen angeben

An die

Vereinigung Deutscher
Elektrizitätswerke VdEW e.V.
z. Hd. Herrn Dr. XXXXX XXX

F r a n k f u r t
Stresemannallee 23

53 Bonn 12, den 30.10.1975

Konery sw 3282
(oder aber Vermittlung 591)

Dienstgebv\$ude: Bonn-Bad Godesberg.
und Bonn, Heuballen 2- 10 (Wohnung)

P-- A--. 7A)

Petr. Endlagerung der Reaktorbetriebsabfv\$11e in ASSE.

A. Je-t Wu New'

chr geehrter Herr r Dr. XXXXXXXXX

No. of clusters	Within cluster variance
6	0.23
5	0.24
4	0.36

↑
Significant increase

34

Example Segmentation

DER BUNDESMINISTER
FOR FORSCHUNG UND TECHNOLOGIE
31'i - 5554 - 7/15
Tel Antwortehreiben Blick Degen Renditezielen angeben

53 Bonn 12, den 30.10.1975
Koneru sw 3282
(oder aber Vermittlung 591)
Dienstgebv/sude: Bonn-Bad Godesberg.
und Bonn, Heuballen 2- 10 (Wohnung)

An die
Vereinigung Deutscher
Elektrizit/swerke VdEW e.V.
z. Hd. Herrn Dr. XXXXX XXX
F r a n'k f u r t
Stresemannallee 23

Petr. Endlagerung der Reaktorbetriebsabfv/§11e in ASSE.
chr geehrter Herr r Dr. XXXXXXXXX

P-- A--. 7A)

A. Je-t Wu New'

Chosen clustering →

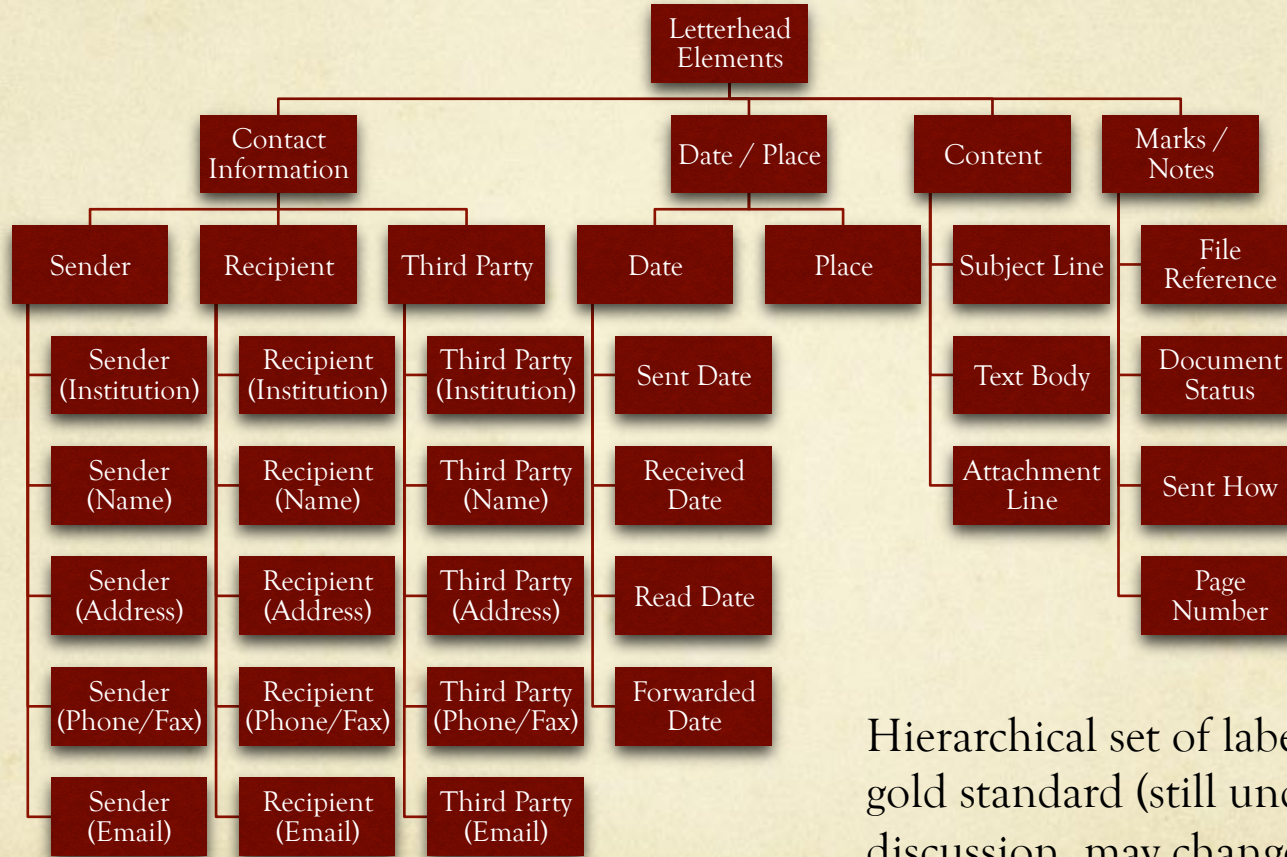
No. of clusters	Within cluster variance
6	0.23
5	0.24
4	0.36

Future Work



- Currently: Classify letterhead segments according to their “semantics” (address, date etc.) using Naïve Bayes and similar
- Layout and content analysis (possibly) not separate
- Next step: Create a gold standard for a small set of letters (ca. 60) and letterhead elements (ca. 1000-1500)
- Evaluate against gold standard

Future Work



Hierarchical set of labels for gold standard (still under discussion, may change)

Future Work



- Currently: Classify letterhead segments according to their “semantics” (address, date etc.) using Naïve Bayes and similar
- Layout and content analysis (possibly) not separate
- Next step: Create a gold standard for a small set of letters (ca. 60) and letterhead elements (ca. 1000-1500)
- Evaluate against gold standard

Acta sunt ser.

Thank you for your attention

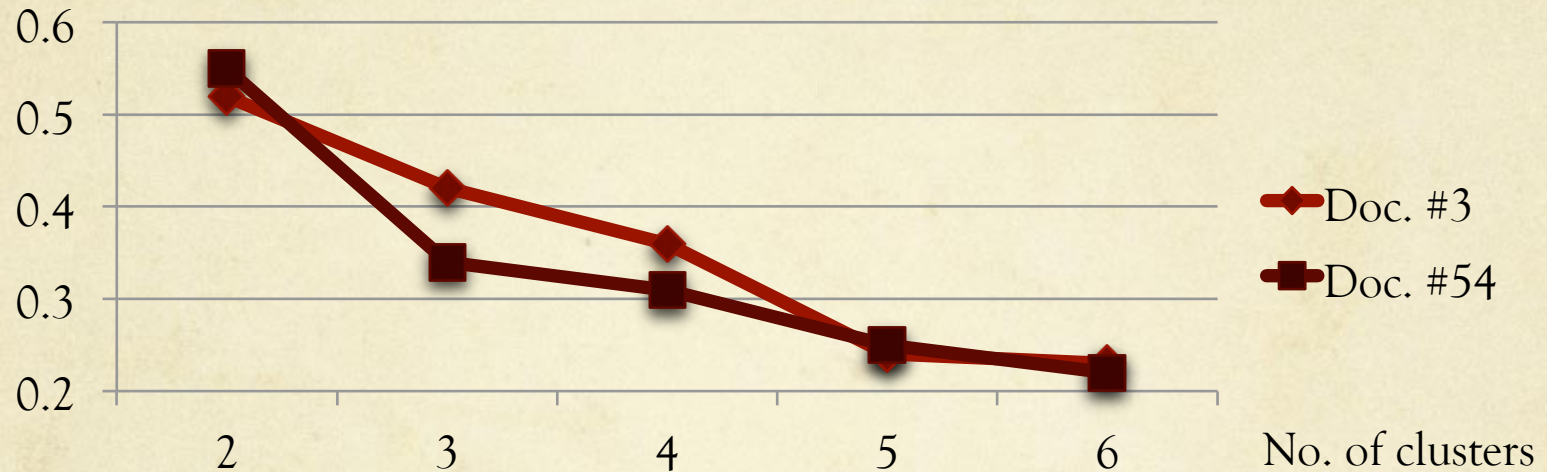
Questions? Comments?



Back-Up / Alternative Slides

Method

Avg. within
cluster
variance



- “Flatten” clustering by identifying significant increase in average within cluster variance (“spot the elbow joint”)
- This indicates that two clusters far away from each other have (erroneously?) been merged

41