



# Nicht immer ist der SOLR schuld!

## Qualitätskontrollen bei Metadaten und Optimierungsmöglichkeiten

Viktor Holzwert, Steffen Illig, Stefan Philipp, **Philipp Rumpf** | Universitätsbibliothek Bamberg

DSpace Praxistreffen

04.–05.04.2024 |  
Mainz

# Die Otto-Friedrich-Universität Bamberg

- seit 1979 (wieder) Universität
- 4 Fakultäten
- ca. 12.000 Studierende
- in der Weltkulturerbestadt Bamberg

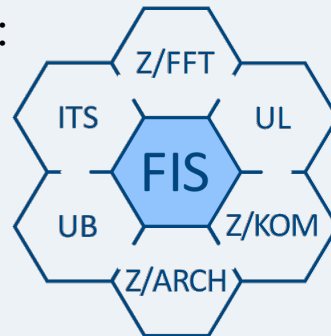


# DSpace in Bamberg – [fis.uni-bamberg.de](https://fis.uni-bamberg.de)

2016	Entscheidung der Universität Bamberg für DSpace-CRIS
2018	Beitritt zum DSpace-Konsortium-Deutschland
2018	Go-Live mit DSpace-CRIS 5.10
2023	Umstieg auf DSpace-CRIS 2023.01.01 (DSpace 7.5)
<b>2024</b>	In Vorbereitung: DSpace-CRIS 2023.02.02 (LTS, DSpace 7.x) Fokus auf: Stabilisierung des Systems

Das FIS ist ein Kooperationsprojekt:

- Forschungsförderung (Z/FFT)
- IT-Service (ITS)
- Universitätsbibliothek
- u.a.



es beinhaltet:


- Bibliografie
- Publikationsserver
- ...

# Problemstellung

Universität Bamberg  
**FIS** 🔍 Anmelden ▾

[Startseite](#) [Publikationen](#) [Forschungsdaten](#) [Projekte](#) [Personen](#) [Einrichtungen](#) [Auszeichnungen](#) [Mein FIS](#) [Hilfe](#)

[Startseite](#) • [Suche](#)


 

🔍 Suche

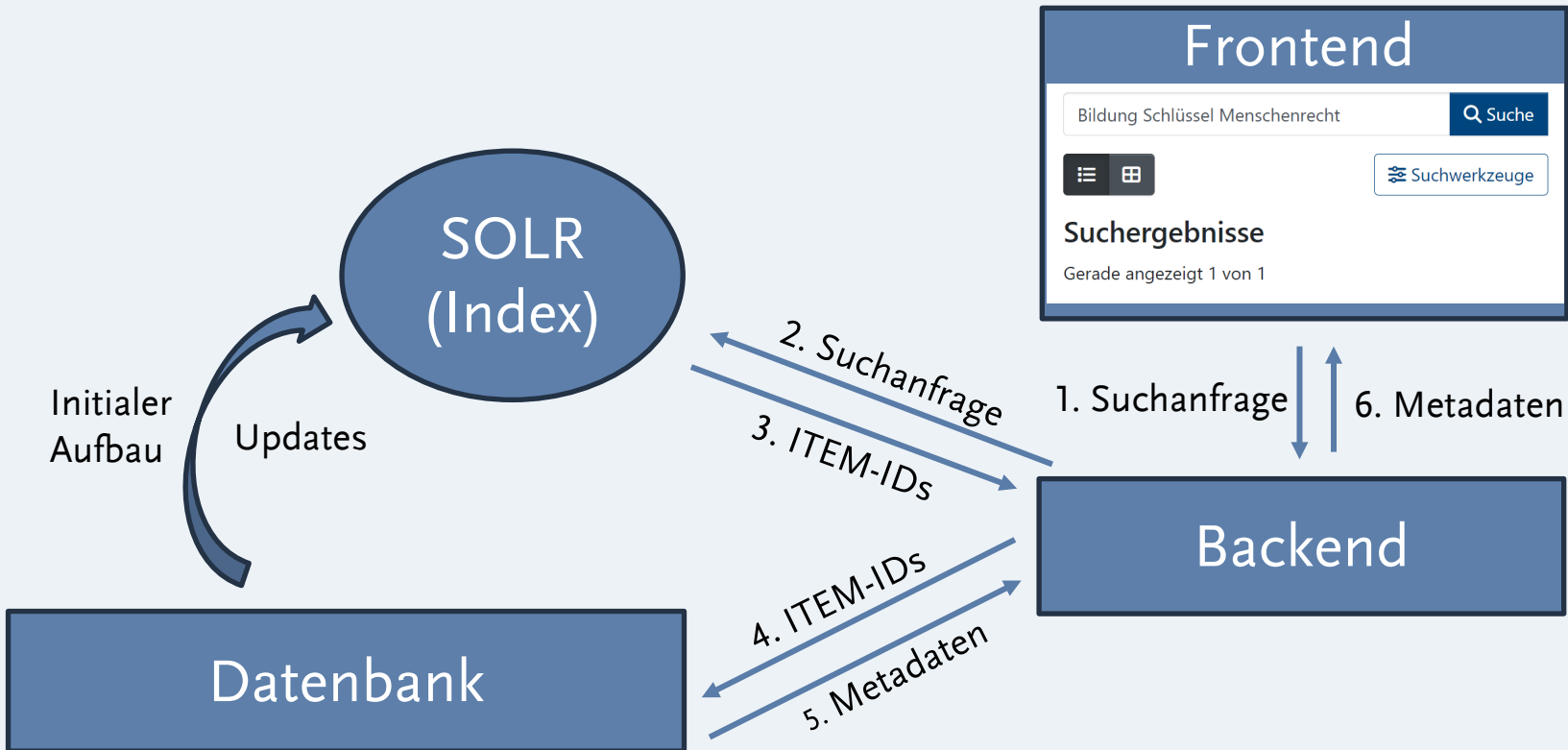
## Suchergebnisse

Ihre Suche führte zu keinem Ergebnis. Versuchen Sie es mit [Anführungszeichen](#)

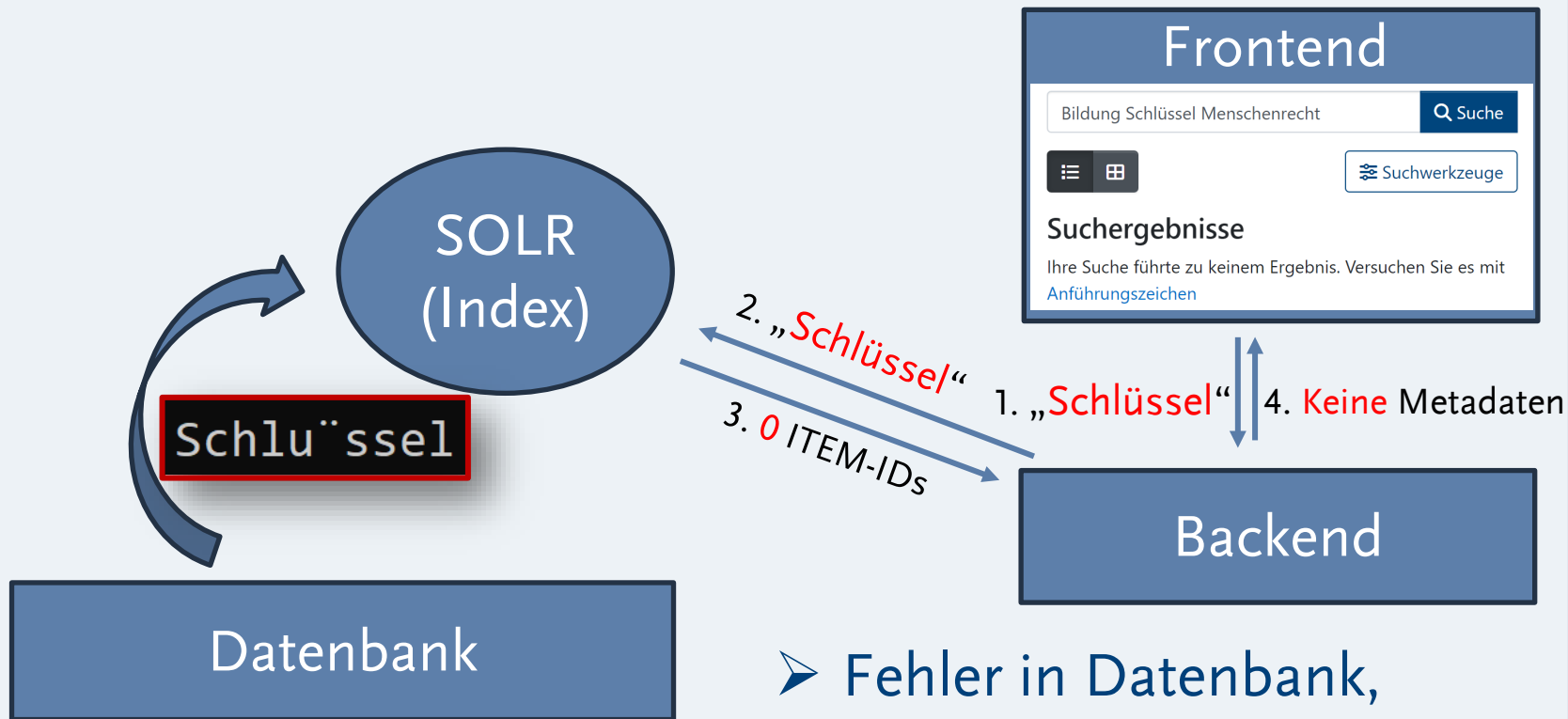
Suche im FIS führt nicht zum gewünschten Treffer, obwohl das ITEM existiert.

Warum? 

# Funktionsweise des SOLR in DS 7 (vereinfacht)



# Funktionsweise des SOLR in DS 7 (vereinfacht)



➤ Fehler in Datenbank, führt zu Fehler im SOLR

# Verschiedene Darstellungen in Unicode (trotz UTF-8 als Standard)

Beispieltitel: Bildung ist der entscheidende Schlüssel

1. Aktuell durch copy-paste „Schlüssel“

→ ü = [75 cc 88] u + Umlaut über vorherigem Buchstaben

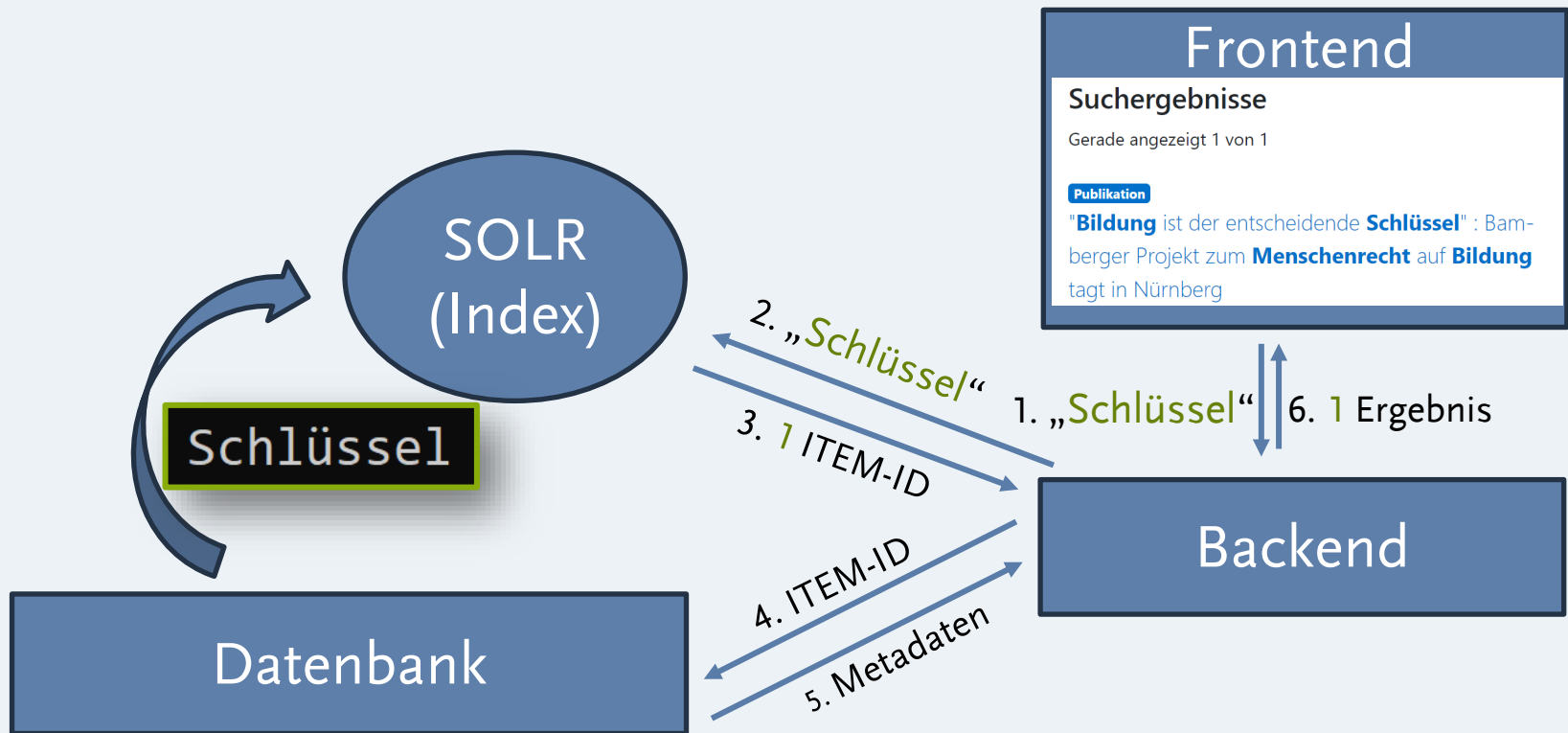
der	entscheidende	Schlüssel
[64 65 72]	[65 6e 74 73 63 68 65 69 64 65 6e 64 65]	[53 63 68 6c 75 cc 88 73 73 65 6c]

2. Eingabe über Tastatur (deutsch) „Schlüssel“ → ü = [c3 bc]

der	entscheidende	Schlüssel
[64 65 72]	[65 6e 74 73 63 68 65 69 64 65 6e 64 65]	[53 63 68 6c c3 bc 73 73 65 6c]

→ Selbes Problem bei ö und ä










# Funktionsweise des SOLR in DS 7 (vereinfacht)














# Weiteres Fehlerbild: Leerzeichen

- Durch falsche Nutzereingabe: Ein oder mehrere Leerzeichen zu viel
- Durch copy-paste: Andere Leerzeichenarten
- Gut wäre ein einzelnes, einheitliches Leerzeichen.

<code>\u0020</code>		space
<code>\u00A0</code>		no-break space
<code>\u2000</code>		en quad
<code>\u2001</code>		em quad
<code>\u2002</code>		en space
<code>\u2003</code>		em space
<code>\u2004</code>		three-per-em space
<code>\u2005</code>		four-per-em space
<code>\u2006</code>		six-per-em space

<code>\u2007</code>		figure space
<code>\u2008</code>		punctuation space
<code>\u2009</code>		thin space
<code>\u200A</code>		hair space
<code>\u202F</code>		narrow no-break space
<code>\u205F</code>		medium mathematical space
<code>\u3000</code>		ideographic space
<code>\u0009</code>		character tabulation
<code>\u1680</code>		ogham space mark



<code>\u0020</code>		space
---------------------	---	-------

# Besonders gemein

Manche Leerzeichenarten wie das geschützte Leerzeichen (non-breaking space) sind auch in der Konsole nicht als falsch zu erkennen:

## Towards the Evaluation of Action Reversibility in STRIPS Using Domain Generators

```
-Towards the Evaluation of Action Reversibility in STRIPS Using Domain Generators  
+Towards the Evaluation of Action Reversibility in STRIPS Using Domain Generators
```

```
Towards the Evaluation of Action Reversibility  
in STRIPS Using Domain Generators
```

# Weiteres Fehlerbild: Breitenlose Leerzeichen (zero width space)

Ansatz neue Lösungen für bestehende testtheoretische Probleme und liefert zugleich konkrete Implikationen für die Praxis der Testkonstruktion und Testauswertung.

Implikationen für die Praxis der Testkonstruktion und Testauswertung.  
Implikationen für die Praxis der Testkonstruktion und Testauswertung.

Implikationen für die Praxis der Testkonstruktion und Testauswertung. `&ZeroWidthSpace;`

# Weiteres Fehlerbild: Breitenlose Leerzeichen (zero width space)

Der SOLR erkennt zwei Worte, wenn das Zeichen in einem Wort vorkommt:  
Auflage → Auf+lage

<code>\u200B</code>	zero width space
<code>\u200C</code>	zero width non-joiner
<code>\u200D</code>	zero width joiner
<code>\u2060</code>	word joiner
<code>\uFEFF</code>	zero width non-breaking space



# Weiteres Fehlerbild: Unterschiedliche Zeilenumbrüche

<code>\u000B</code>	line tabulation
<code>\u000C</code>	form feed
<code>\u0085</code>	next line
<code>\u2028</code>	line separator
<code>\u2029</code>	paragraph separator

Ausreißer

`\u000A`  
line feed

Linux  
macOS

`\u000D\u000A`  
carriage return + line feed

Windows

# Weiteres Fehlerbild: Unterschiedliche Zeilenumbrüche

Beispiele:

for focused disciplinary and interdisciplinary research projects.

Our approach presents a novel concept for data federation in

```
interdisciplinary research projects.\nOur approach pr  
interdisciplinary research projects.\r\nOur approach
```

are responsible for reducing obsessive-compulsive symptoma-  
tology. After an initial exploration of the literature, I hypothesi-

```
symptomatology. \nAfter an initial  
symptomatology.\r\nAfter an initial
```

# Wir analysieren unsere Datenbank ...

... mit einem selbstentwickelten Skript für **DSpace 7**

## ➤ DSpace Pull Request #9427

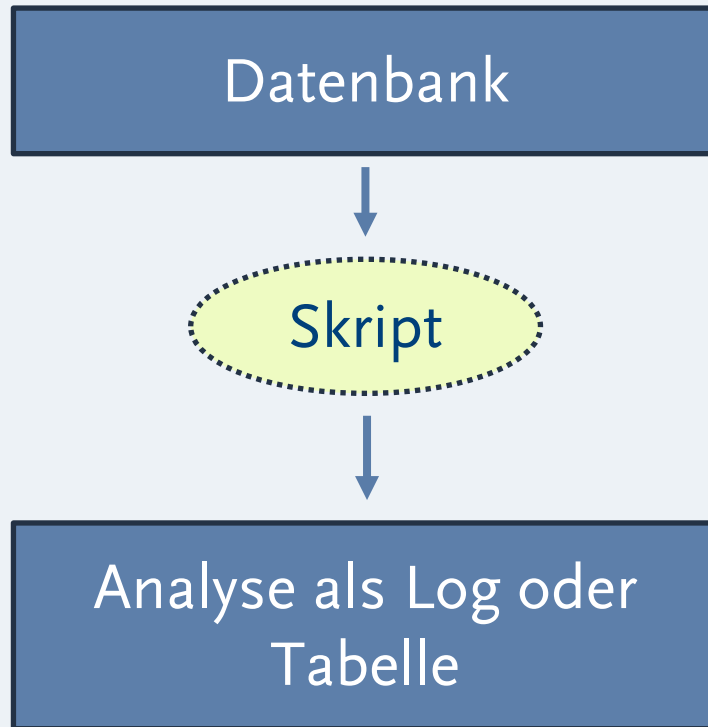
Das Skript verwendet die Java Klasse „Enhancer“ (nicht die Enhancer Funktionalität von DS-CRIS) und kontrolliert in unserer Datenbank folgende Felder:

- dc.title
- dc.title.alternative
- dc.description.abstract
- (weitere Felder sind möglich)

Das Skript kann:

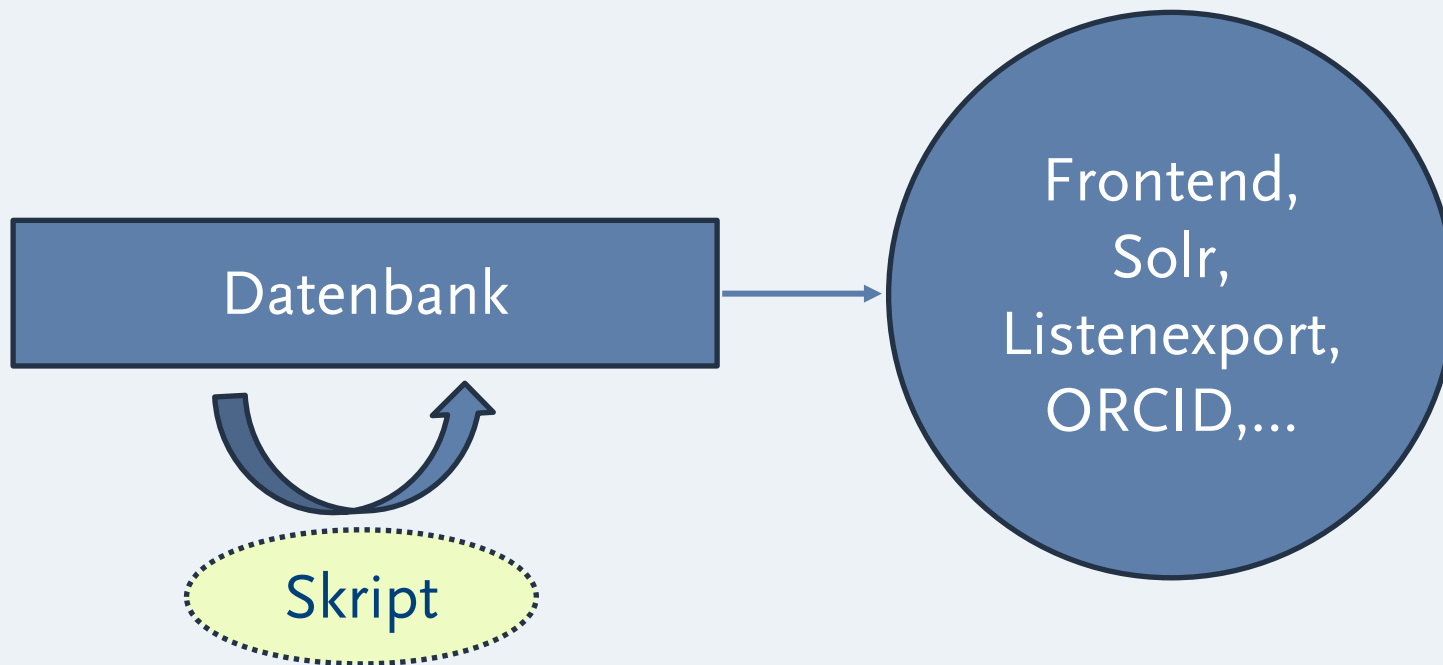
- Eine Ergebnisliste anzeigen ==> Für manuelle Prüfung
- Fehler automatisch bereinigen ==> wenn gewünscht

# Das Skript analysiert ...





... und korrigiert auf Wunsch (optional)



# Zum Beispiel: Unterschiedliche Zeilenumbrüche

➤ Zeilenumbrüche werden vom Skript in 2 Schritten korrigiert:

1. Alle Zeichen, die Zeilenumbrüche markieren, werden durch Zeilenvorschub (englisch: line feed) ersetzt.

2. Zeilenvorschübe werden zu Zeilenumbrüchen (Windows Standard) kombiniert

<code>\u000B</code>	line tabulation
<code>\u000C</code>	form feed
<code>\u0085</code>	next line
<code>\u2028</code>	line separator
<code>\u2029</code>	paragraph separator

**Ausreißer**



<code>\u000A</code>
line feed

**Linux**  
**macOS**



<code>\u000D\u000A</code>
carriage return + line feed

**Windows**

# Lessons learned: Fehler in unserer Datenbank

Datensätze mit:	dc.title	dc.title.alternative	dc.description.abstract
Gesamtzahl der Datensätze	59.836	986	5.545
Betroffen	1.126	9	738
Zeilenvorschübe			52
Spezielle Zeilenumbrüche			2
Leerzeichen vor dem Zeilenumbruch			417
Multiple Leerzeichen	1.075	9	250
Spezielle Leerzeichen	43		118
Umlaute	11		18

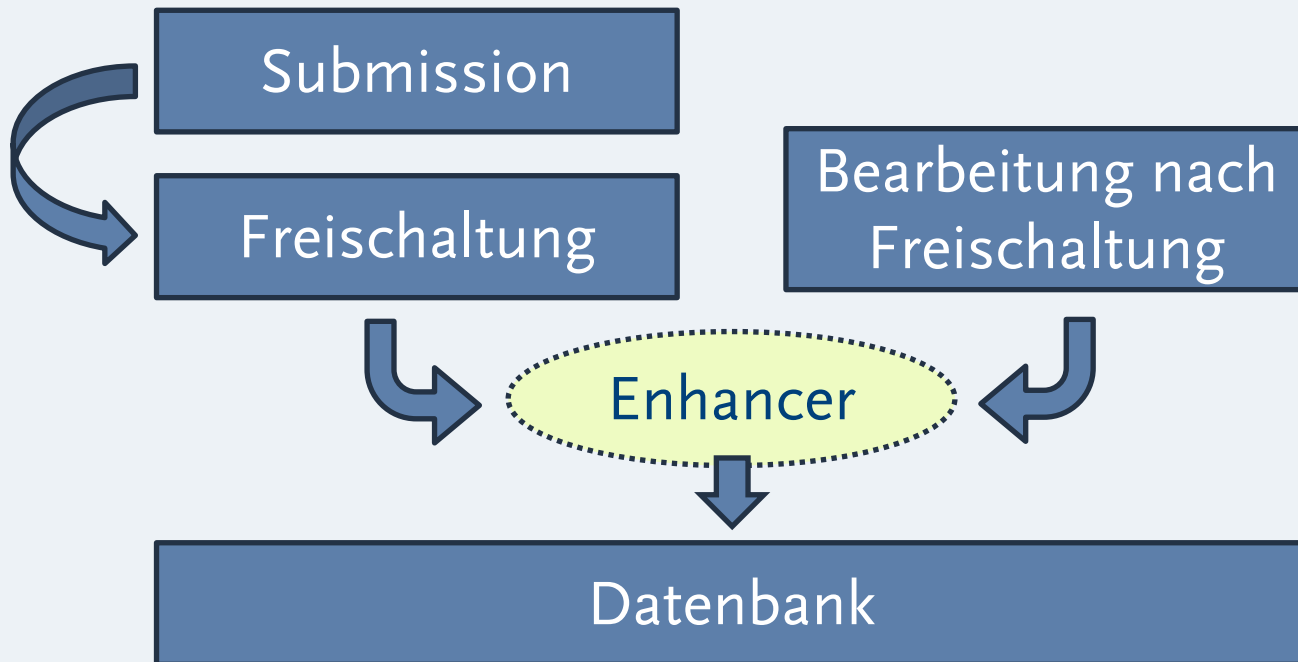
- Das Bamberger FIS ist Bibliografie und Publikationsserver.
- Bei Einträgen mit Dateien (mit Abstract) häufen sich die Fehler. Größeres Problem bei Publikationsservern!

# Ist mit dem Skript das Problem gelöst?

- Die Fehler in unserer Datenbank sind bereinigt. Wir zeigen die Metadaten korrekt an und geben diese korrekt weiter.
  - Es werden aber weiterhin neue Daten aus unterschiedlichen (fehlerbehafteten) Quellen per copy-paste eingegeben.
- Das Problem kehrt wieder. 😞

**Lösung: Wir nutzen die DSpace-CRIS-Funktionalität des „Enhancers“!**

# Der Enhancer korrigiert (anders als das Skript) fortlaufend ab Freischaltung



[DSpace-CRIS Pull Request #444](#)

# Ist ein Enhancer für DSpace 7 denkbar?

Anders als DSpace-CRIS kann DSpace keine Enhancer ausführen.  
Es braucht dazu einen neuen Listener.

Dieser Listener würde unseren Enhancer nachnutzen.

Freiwillige gesucht:

- ➔ Issue eröffnen
- ➔ PR schreiben
- ➔ Feedback der Community erbitten/einarbeiten.



Vielen Dank für die Aufmerksamkeit!  
Sind noch Fragen offen?

