

Wie indexieren Google & Co 13 Millionen Seiten?

Inetbib-Tagung Bonn, 05.11.2004

Florian Seiffert, HBZ



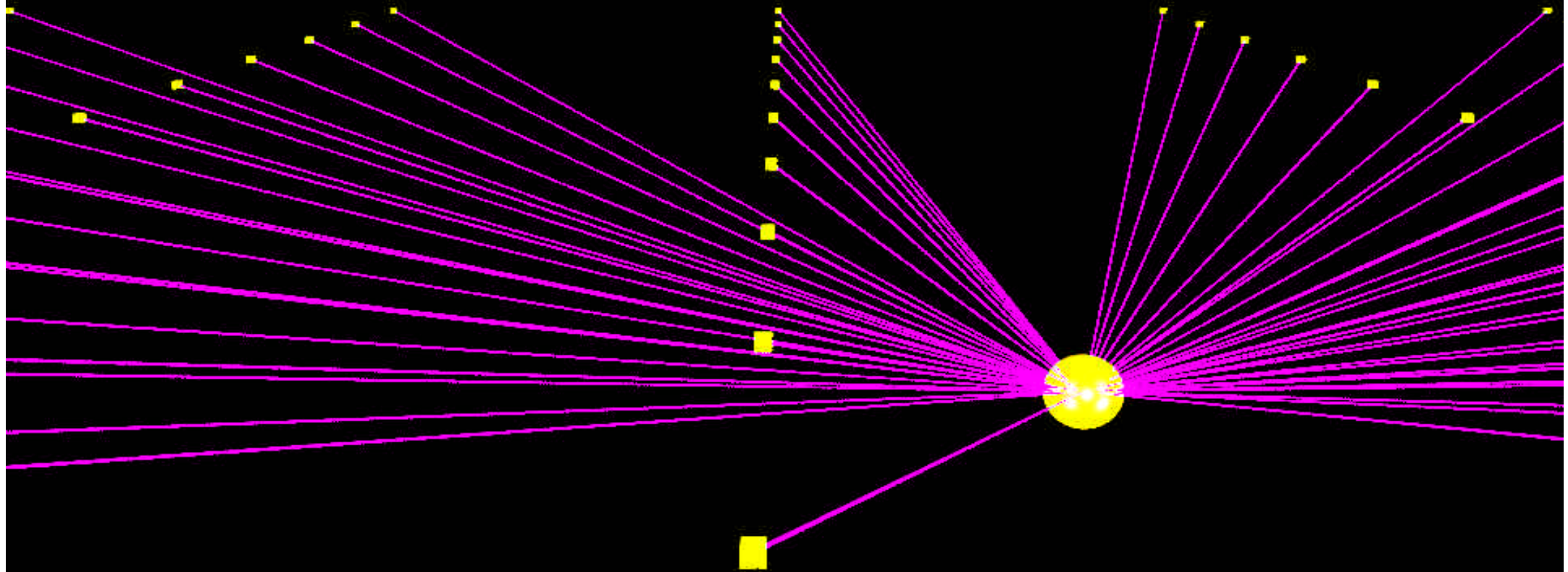
Was haben wir denn heute so vor ...

- Visualisierung
- Google
 - Indexierungsstrategie
 - Indexierungsleistung
 - Seitenzahlbegrenzung pro Server
 - Dauer bis zur Findbarkeit in der Suche
 - "Tiefe" der Indexierung
 - Dateigrößen
 - Einfluss auf die Besuchshäufigkeit der Bots
- fast, msn, neofonie
- Was lernen wir daraus

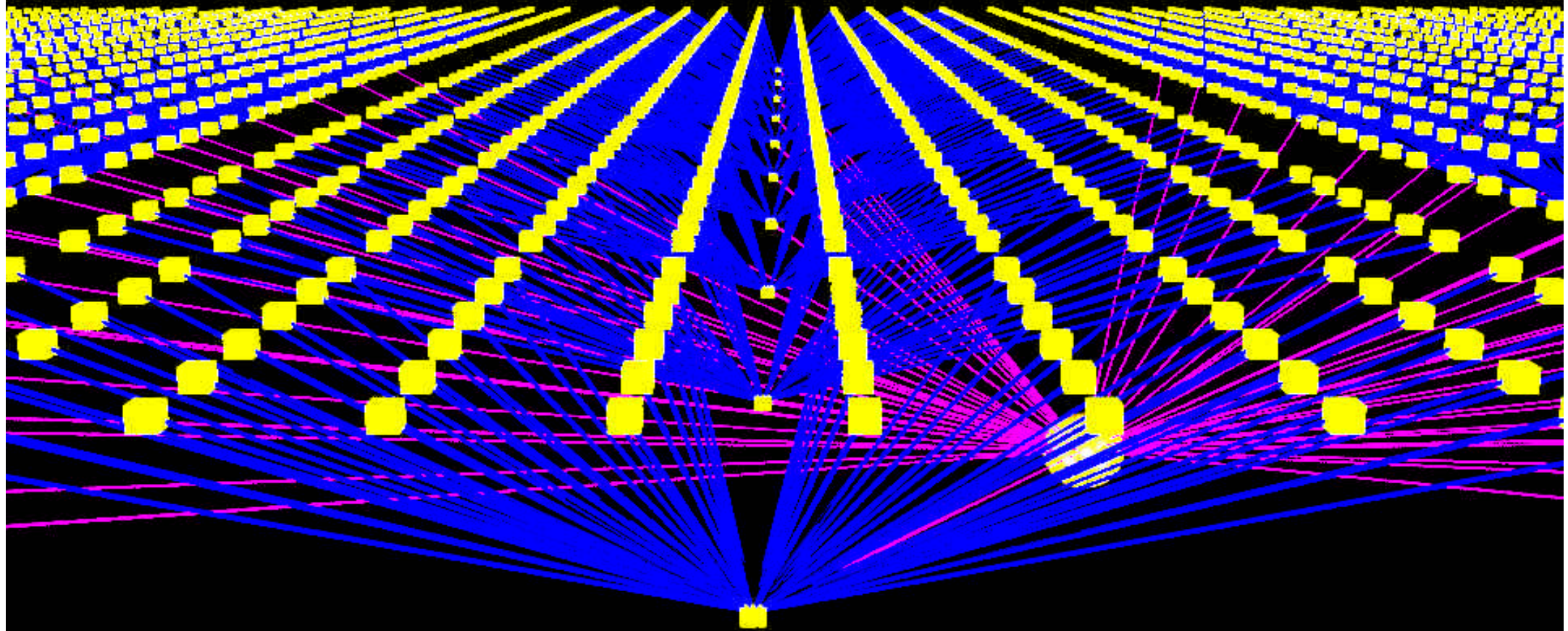
Visualisierung: root



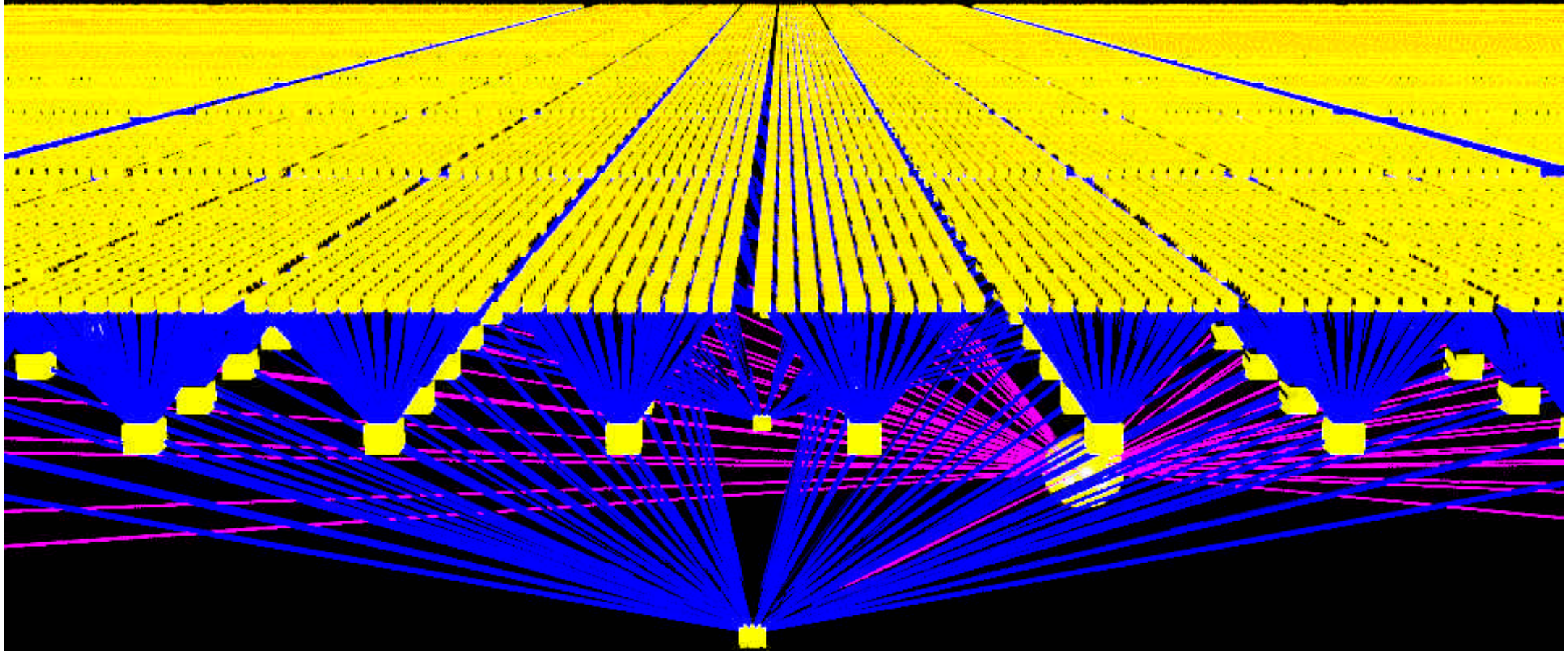
Visualisierung: 60 Seiten



Visualisierung: $60 \times 60 = 3600$ Seiten

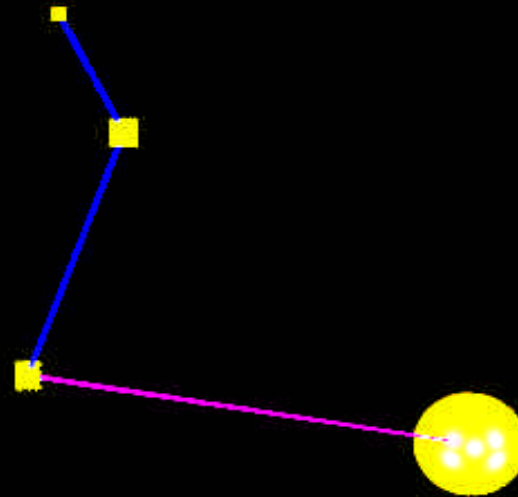


Visualisierung: $60 \times 60 \times 60 = 216.000$ Zellen

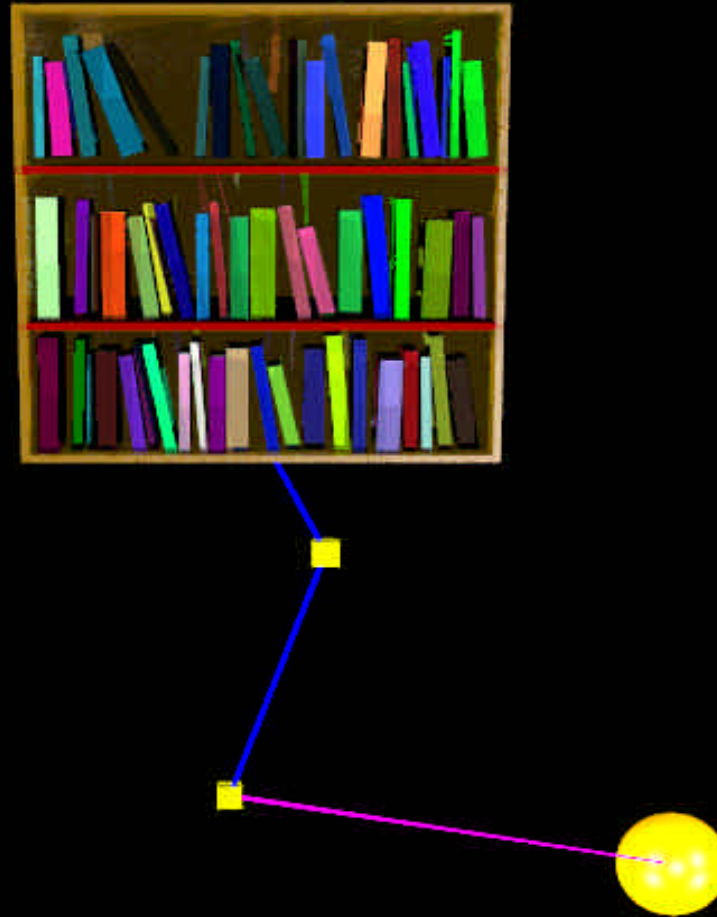


Florian Seiffert, HBZ

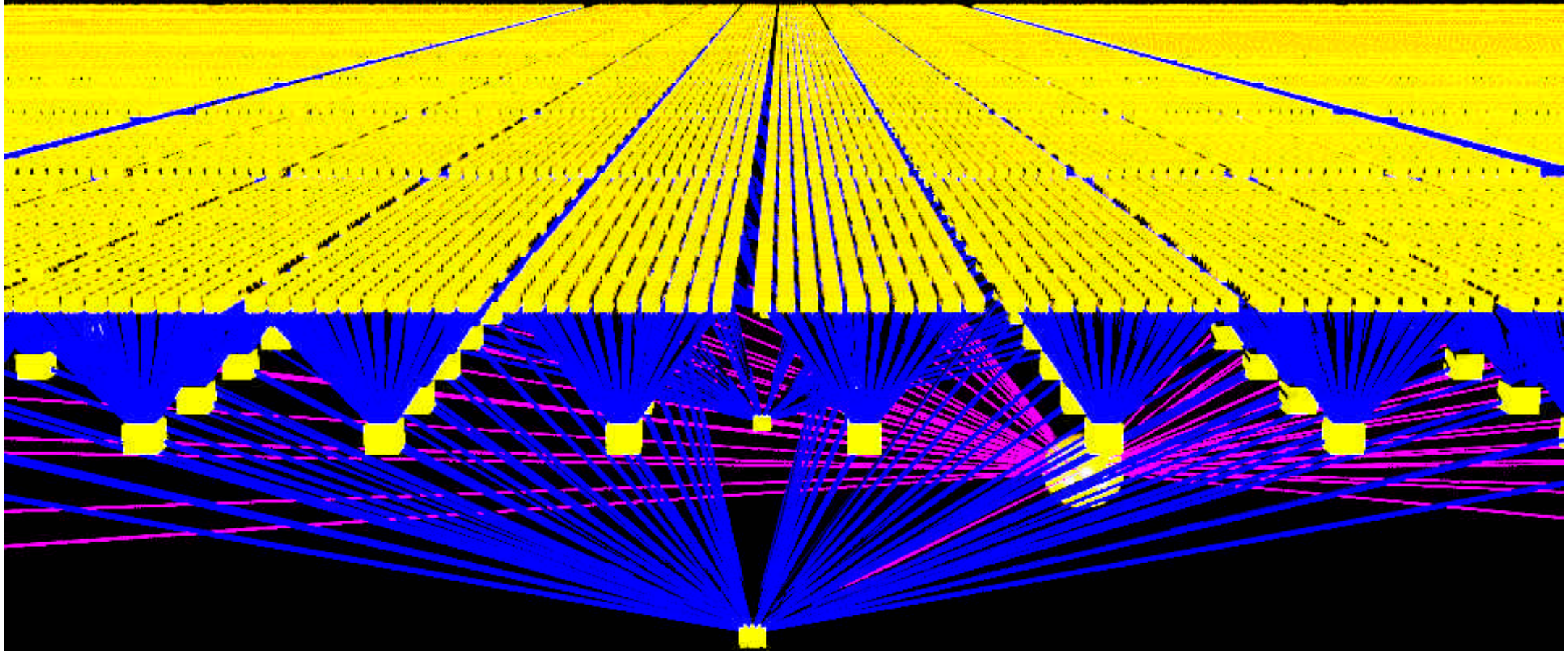
Visualisierung: Eine Zelle



Visualisierung: 60 Buecher pro Zelle -> 12.96 Mio

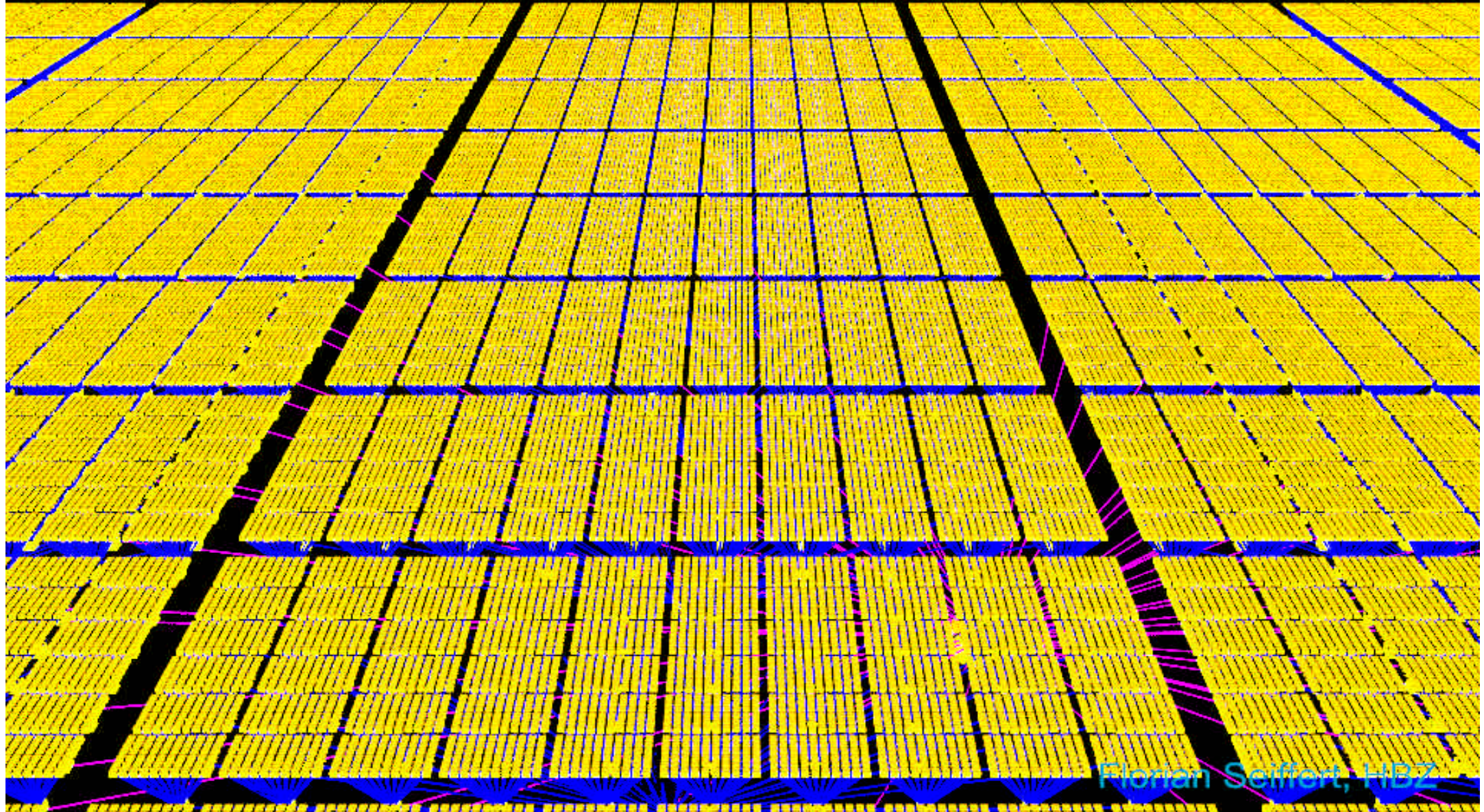


Visualisierung: $60 \times 60 \times 60 = 216.000$ Zellen

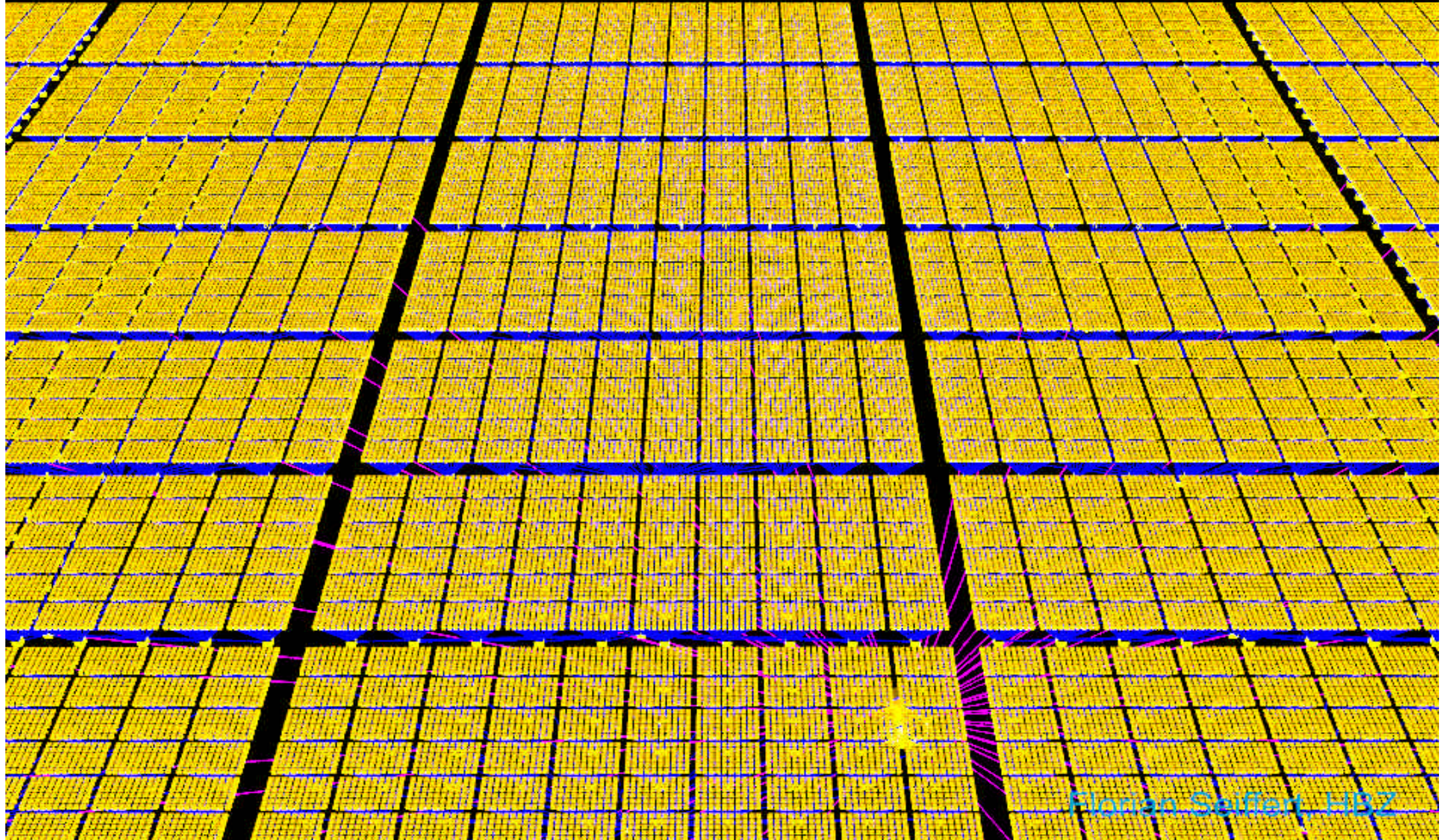


Florian Seiffert, HBZ

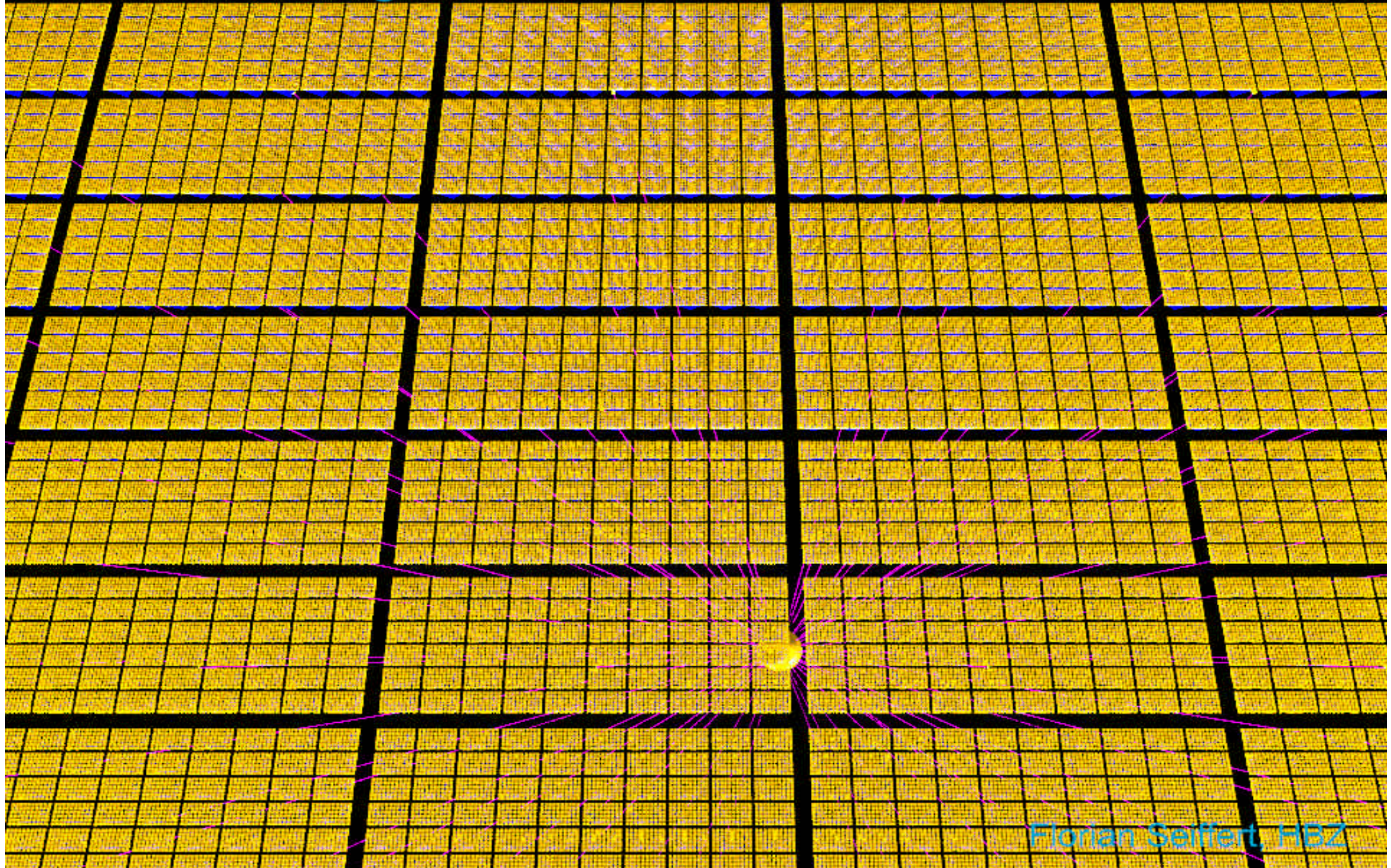
Visualisierung: Wir drehen die Konstruktion



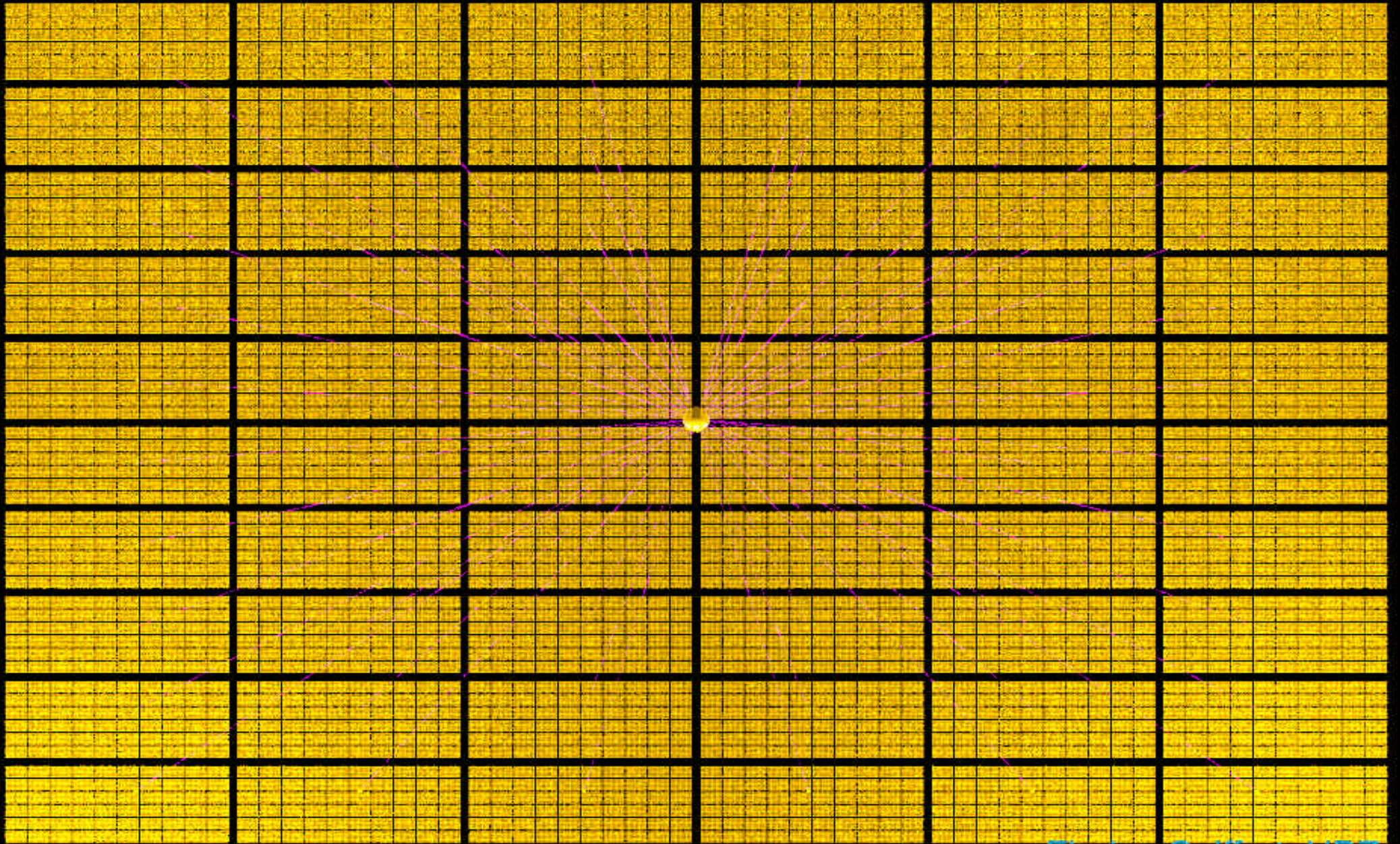
Visualisierung: Wir drehen die Konstruktion



Visualisierung: Wir drehen die Konstruktion



Visualisierung: Fertig

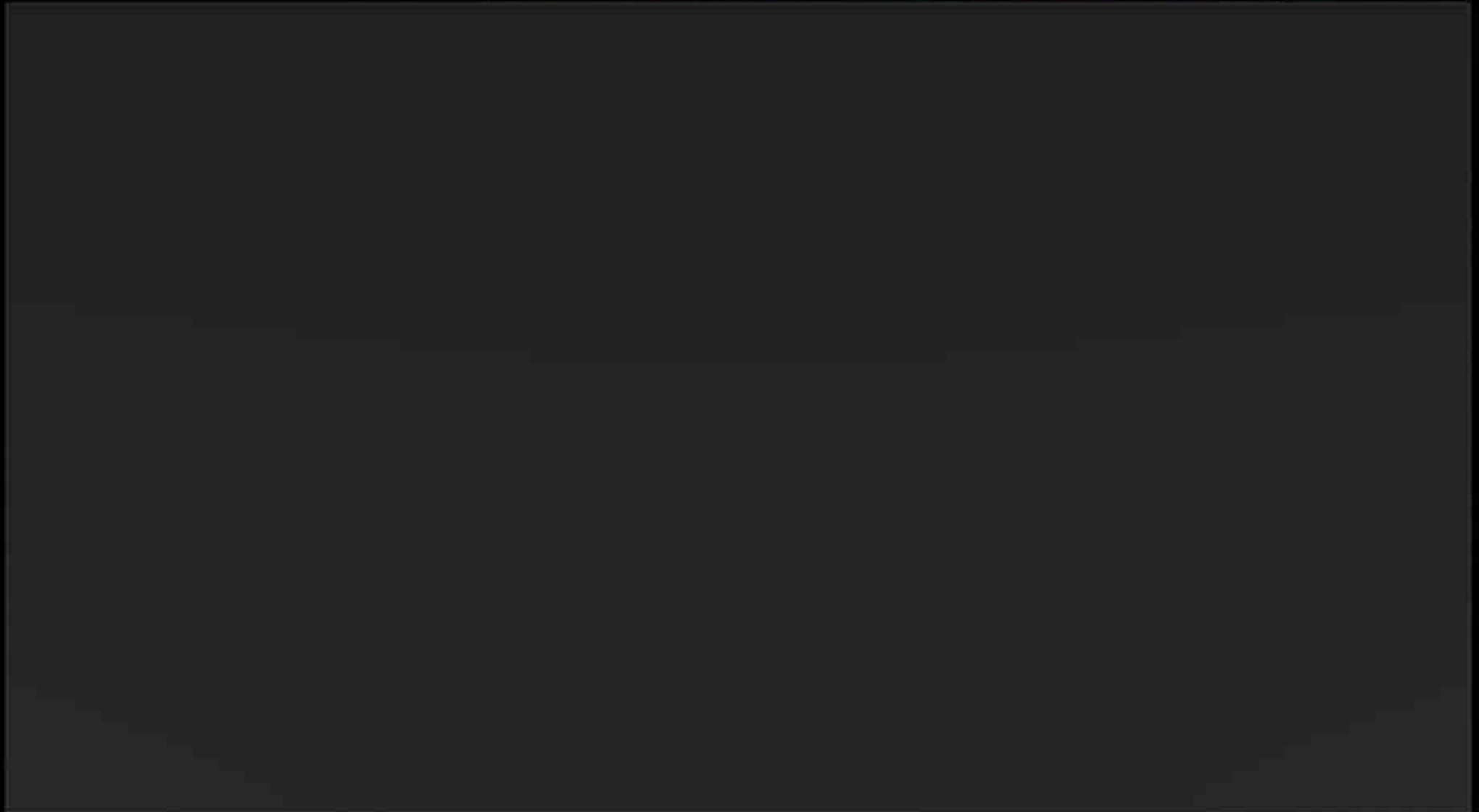
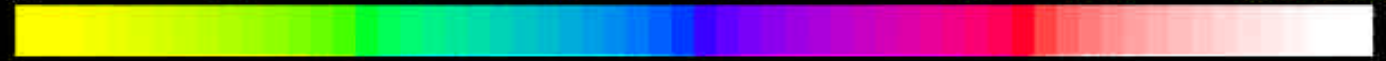


Visualisierung: Abstände entfernt



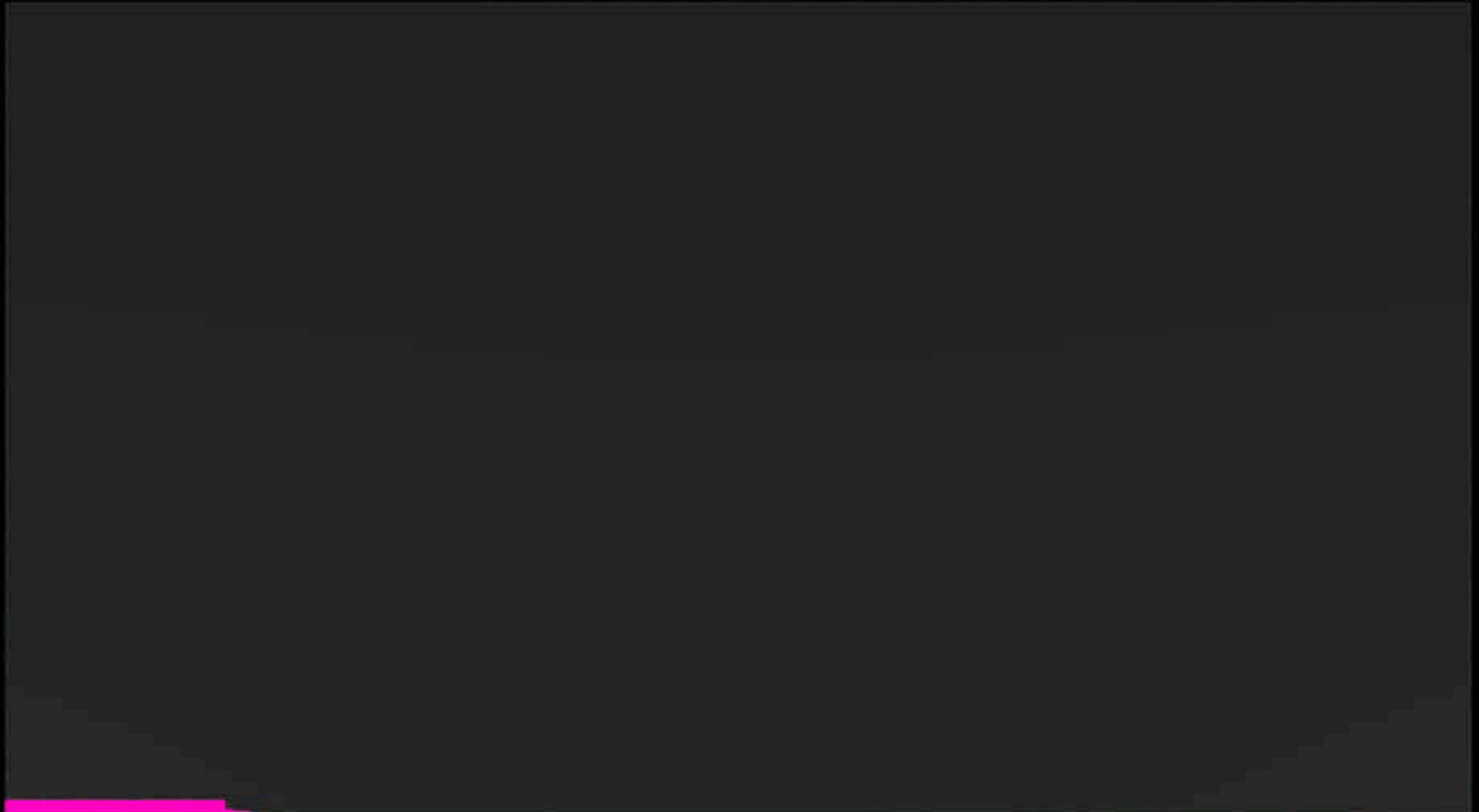
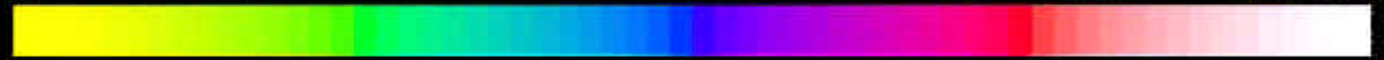
Visualisierung: Monitor auf vbR

Zahl der Treffer pro Zelle: 1 2 5 10 20 60 100 200 300



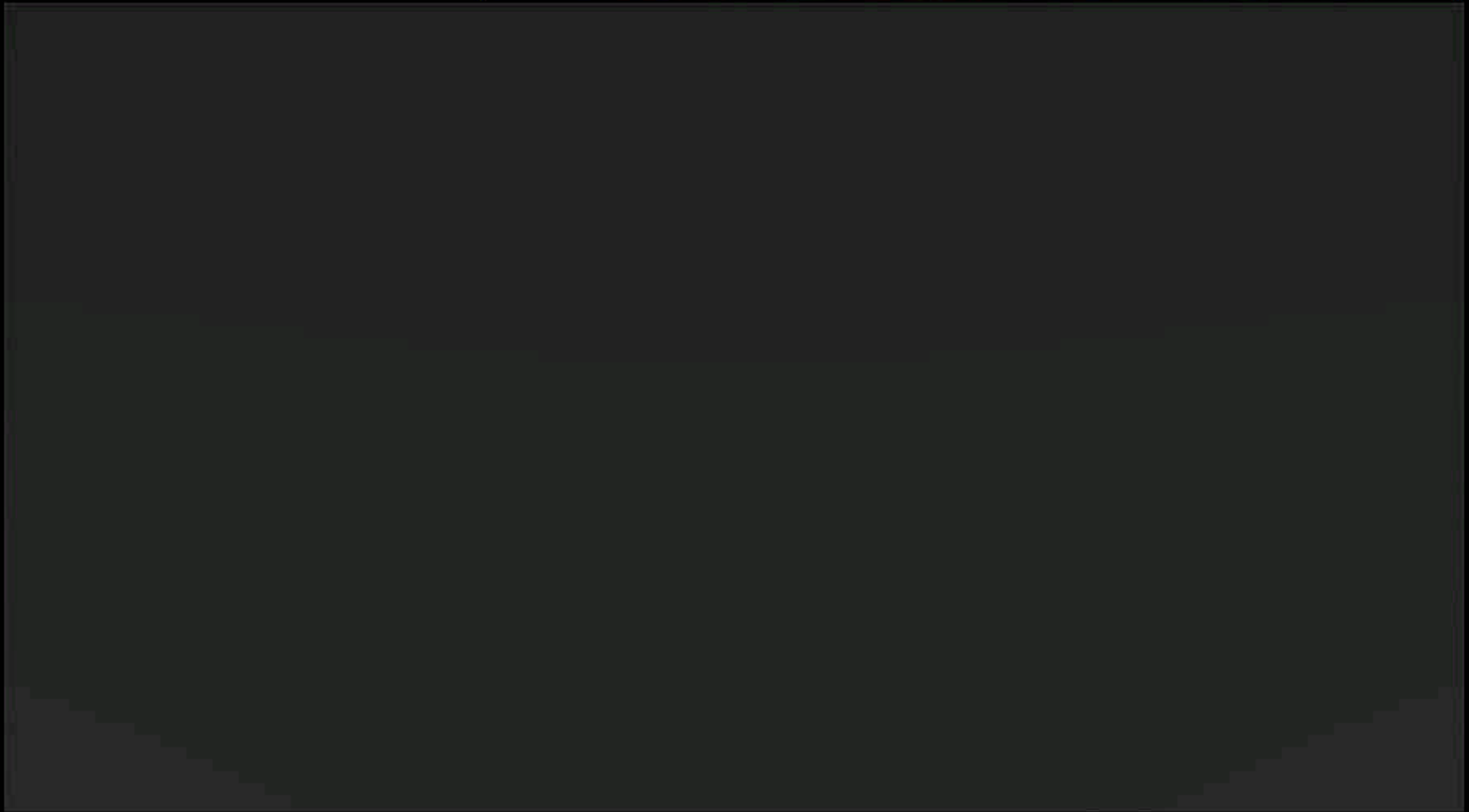
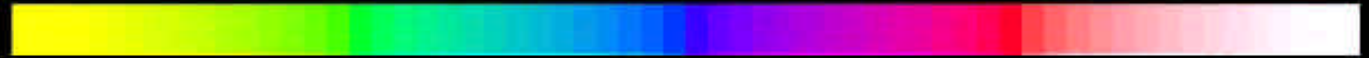
Simulation: linear

Zahl der Treffer pro Zelle: 1 2 5 10 20 60 100 200 300



Googlebot vom 31/Jul/2002

Zahl der Treffer pro Zelle: 1 2 5 10 20 60 100 200 300

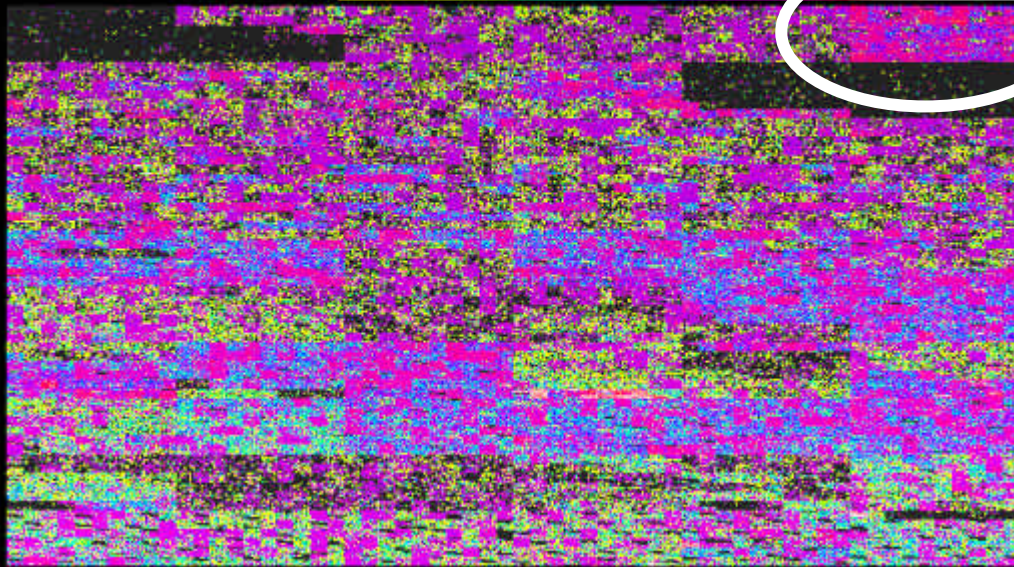


Was haben wir gesehen?

- Keine systematische Indexierung
 - nix: Seite 1, Seite 2,, Seite 13 Mio
 - statistische Verfolgung der Links
- Erklärungsversuche
 - unbenutzte Ids
 - Schlagwoerter

von Google indexiert

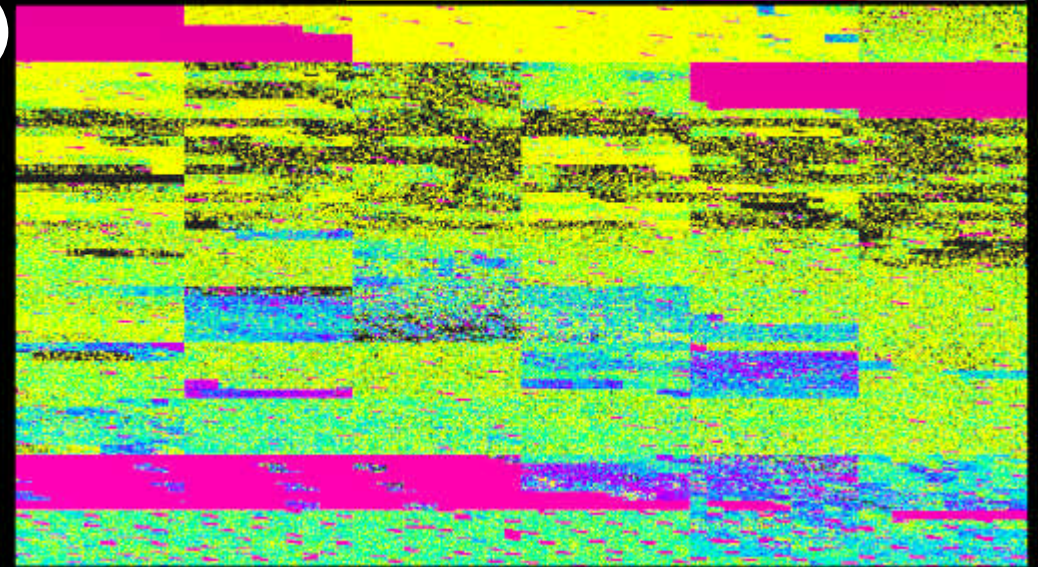
Zahl der Treffer pro Zelle: 1 2 5 10 20 60 100 200 300



Florian Seiffert, HBZ

unbenutzte IDs

Zahl der Treffer pro Zelle: 1 2 5 10 20 60 100 200 300

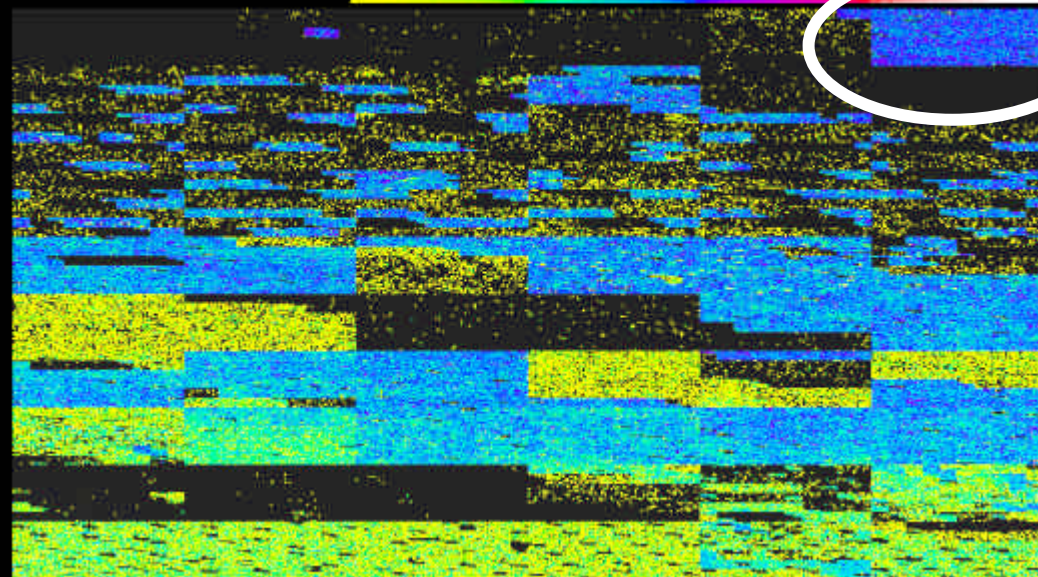


Florian Seiffert, HBZ

- unbenutzte Ids
- Schlagwoerter

Titel mit Schlagworten

Zahl der Treffer pro Zelle: 1 2 5 10 20 60 100 200 300



Florian Seiffert, HBZ

Schlagworte ein Vorteil?

- **Ja!**

- **Google-Bots indexieren Seiten**

- mit Schlagworten: 686.561 (18.2%)

- ohne Schlagworte: 3.090.095 (81.8%)

- **Suchende ueber Google finden Seiten**

- mit Schlagworten: 154.375 (41.7%)

- ohne Schlagworte: 216.072 (58.3%)

- **9.7% der Titel im Verbund sind verschlagwortet**

Schlagworte ein Vorteil?

- **Ja!**

- **Google-Bots indexieren Seiten**

- mit Schlagworten: 686.561 (18.2%)

- ohne Schlagworte: 3.090.095 (81.8%)

18.2%

- **Suchende ueber Google finden Seiten**

- mit Schlagworten: 154.375 (41.7%)

- ohne Schlagworte: 216.072 (58.3%)

41.7%

- **9.7% der Titel im Verbund**

gegen 9.7%

Was haben wir denn heute so vor ...

- Visualisierung
- Google
 - Indexierungsstrategie
 - Indexierungsleistung
 - Seitenzahlbegrenzung pro Server
 - Dauer bis zur Findbarkeit in der Suche
 - "Tiefe" der Indexierung
 - Dateigrößen
 - Einfluss auf die Besuchshäufigkeit der Bots
- fast, msn, neofonie
- Was lernen wir daraus

Was haben wir denn heute so vor ...

- Visualisierung ✓
- Google ✓
 - Indexierungsstrategie ✓
 - Indexierungsleistung
 - Seitenzahlbegrenzung pro Server
 - Dauer bis zur Findbarkeit in der Suche
 - "Tiefe" der Indexierung
 - Dateigrößen
 - Einfluss auf die Besuchshäufigkeit der Bots
- fast, msn, neofonie
- Was lernen wir daraus

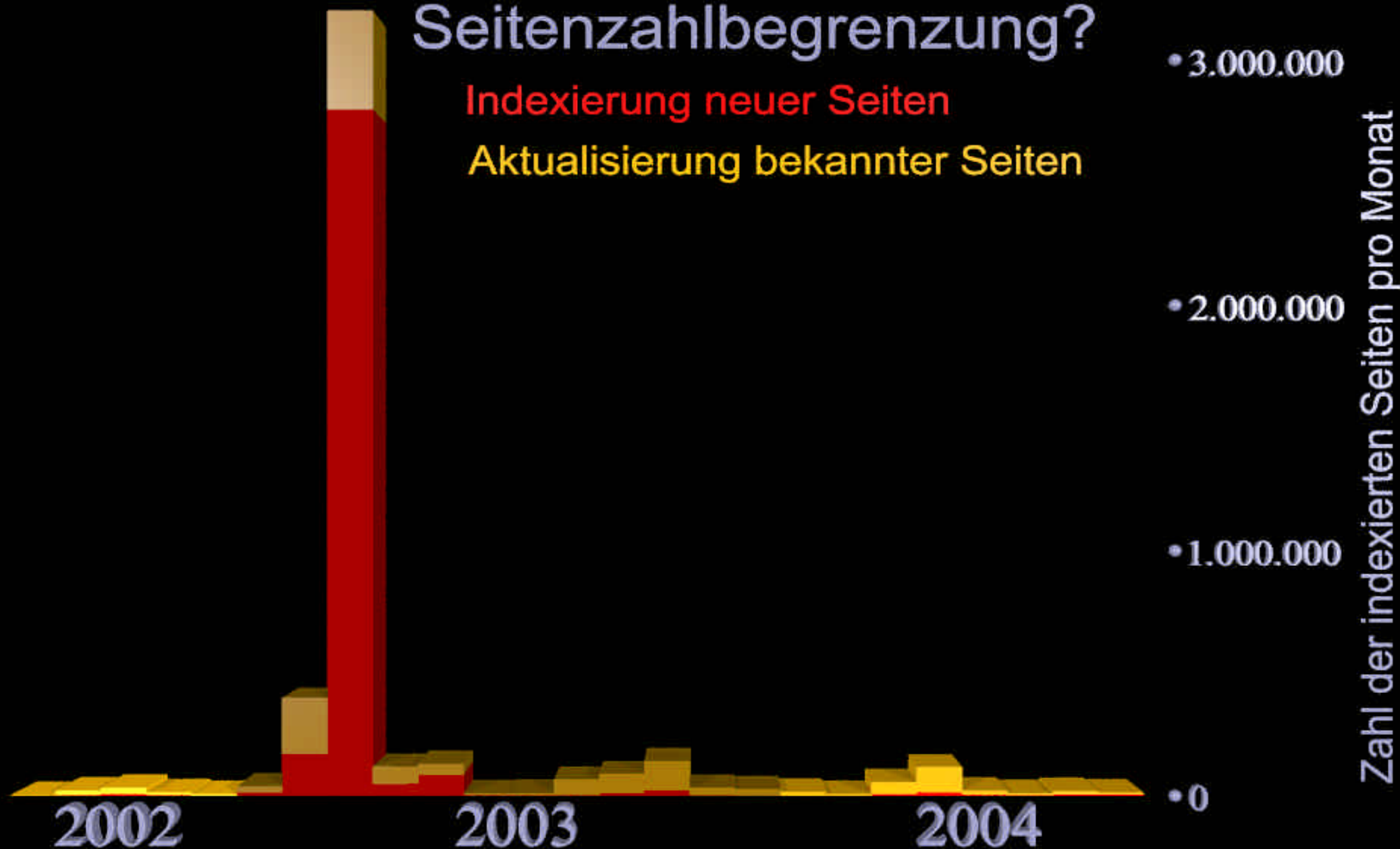
Wieviele Seiten?

- wurden von Google im virt. Buecherregal "abgegrast"?
 - insgesamt 3.776.656 (27.3%)
- kennt Google im virt. Buecherregal?
 - insgesamt 3.220.818 (23.3%)
- wurden durch Google-Suchen gefunden?
 - insgesamt 370.447 (2.7%)

Seitenzahlbegrenzung?

Indexierung neuer Seiten

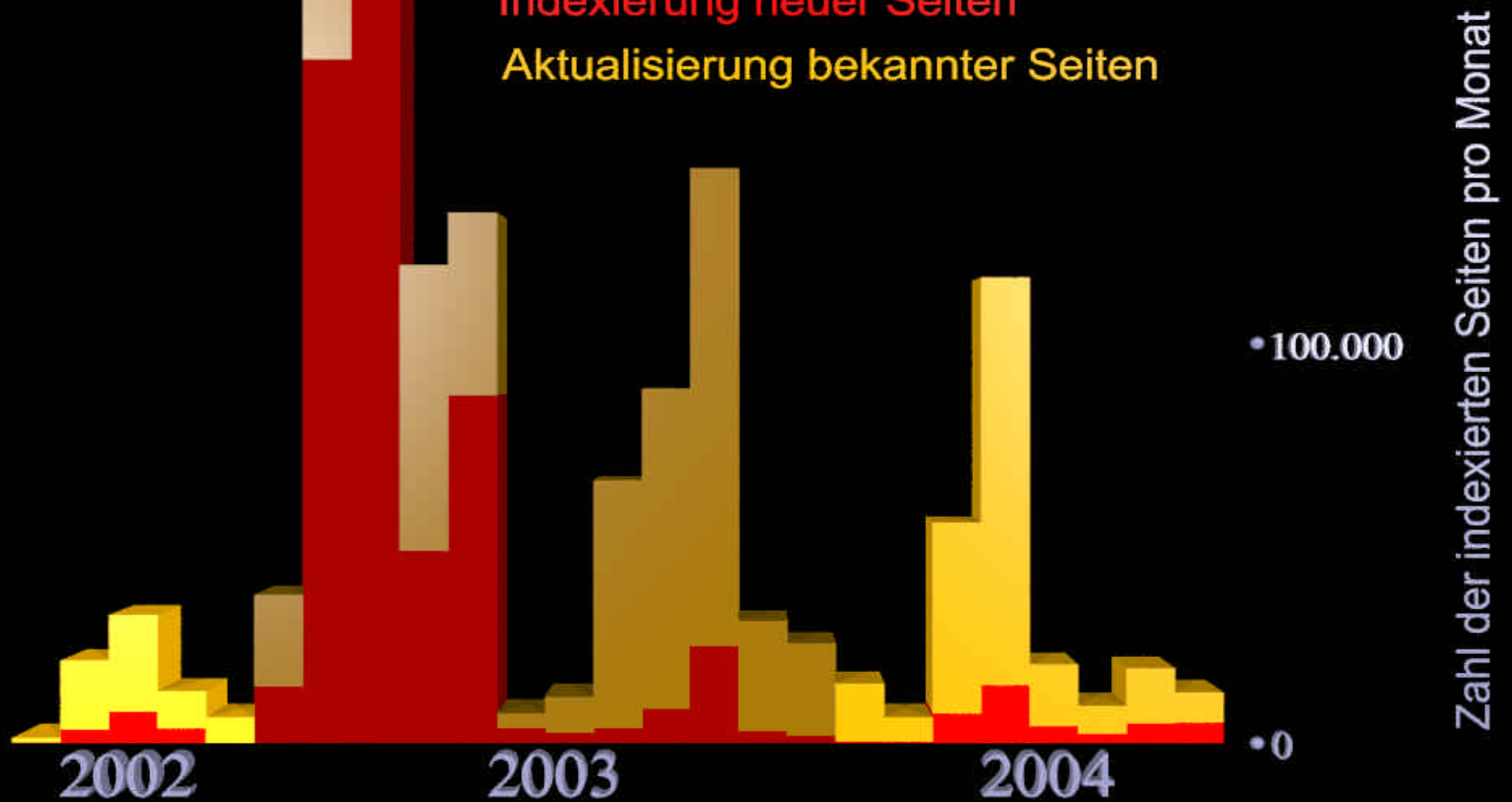
Aktualisierung bekannter Seiten



Seitenzahlbegrenzung?

Indexierung neuer Seiten

Aktualisierung bekannter Seiten



Seitenzahlbegrenzung?

Indexierung neuer Seiten

Aktualisierung bekannter Seiten

79.5 %
20.5 %

25.7 %
74.3 %

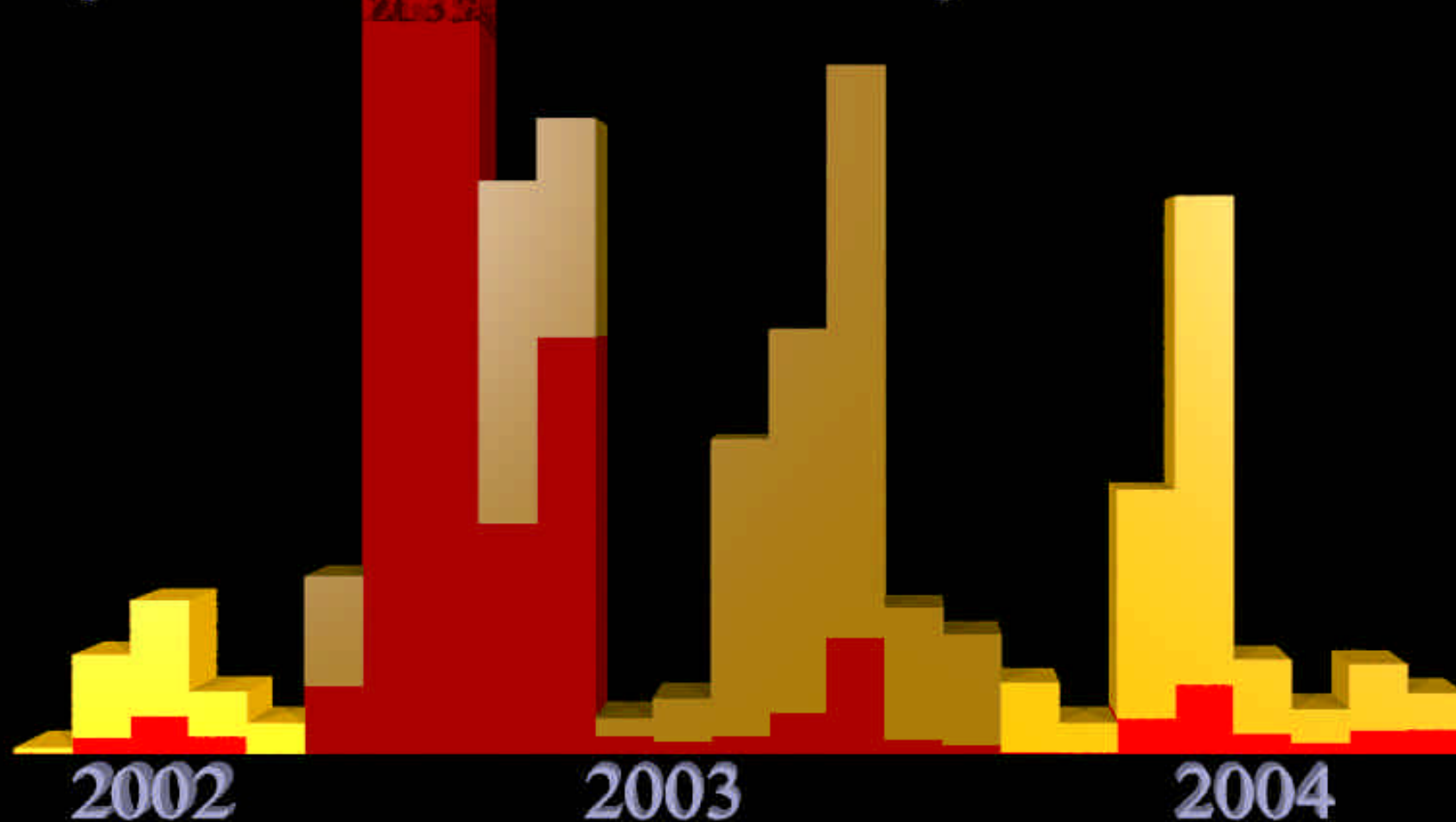
84.5 %
15.5 %



Zahl der indexierten Seiten pro Monat

virt. Buecherregal: Eine Hoechstzahl von Seiten die Google pro Server indiziert ist nicht erkennbar.

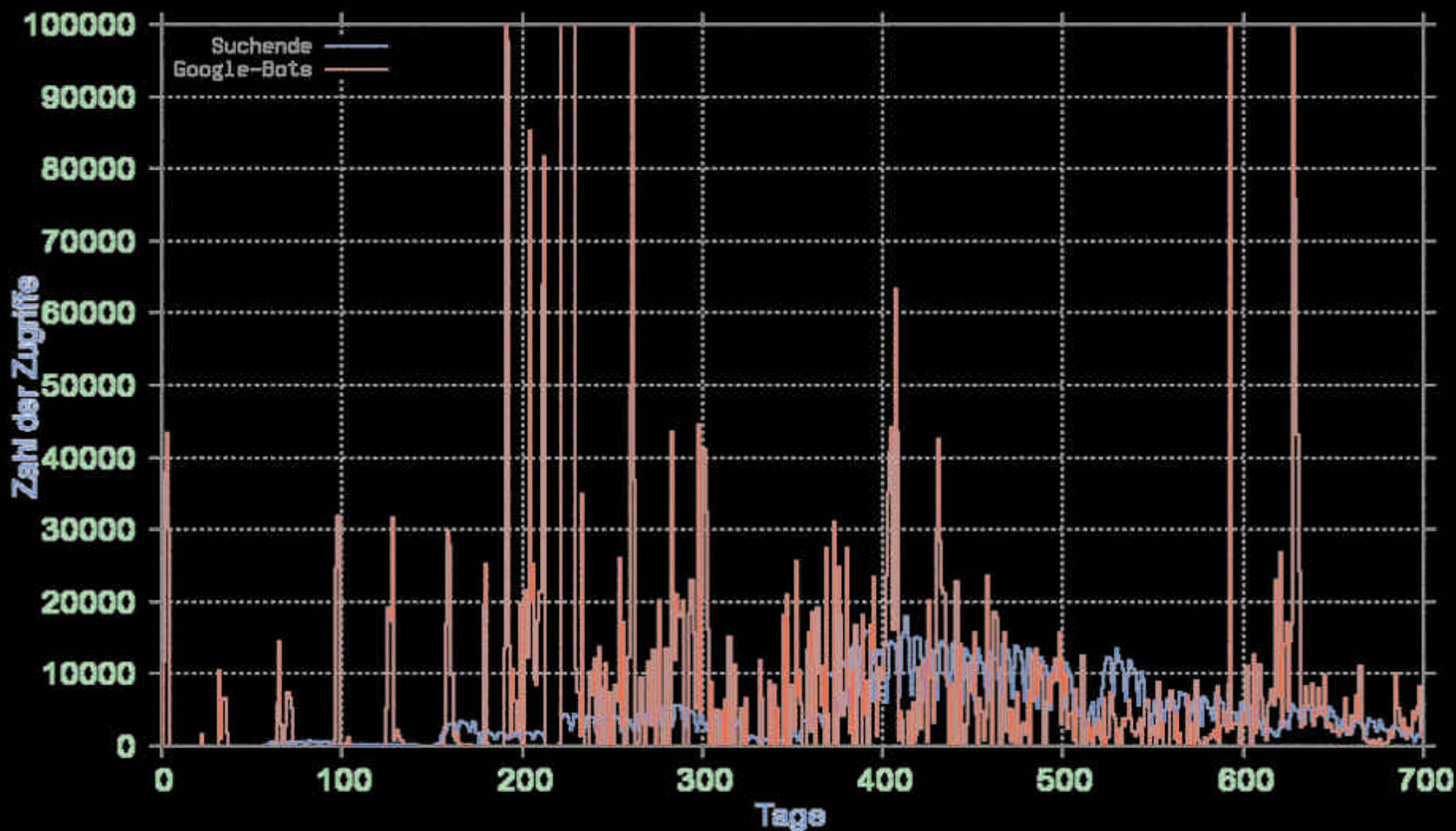
Google kennt z.Zt. 3.2 Mio Seiten im virt. Buecherregal



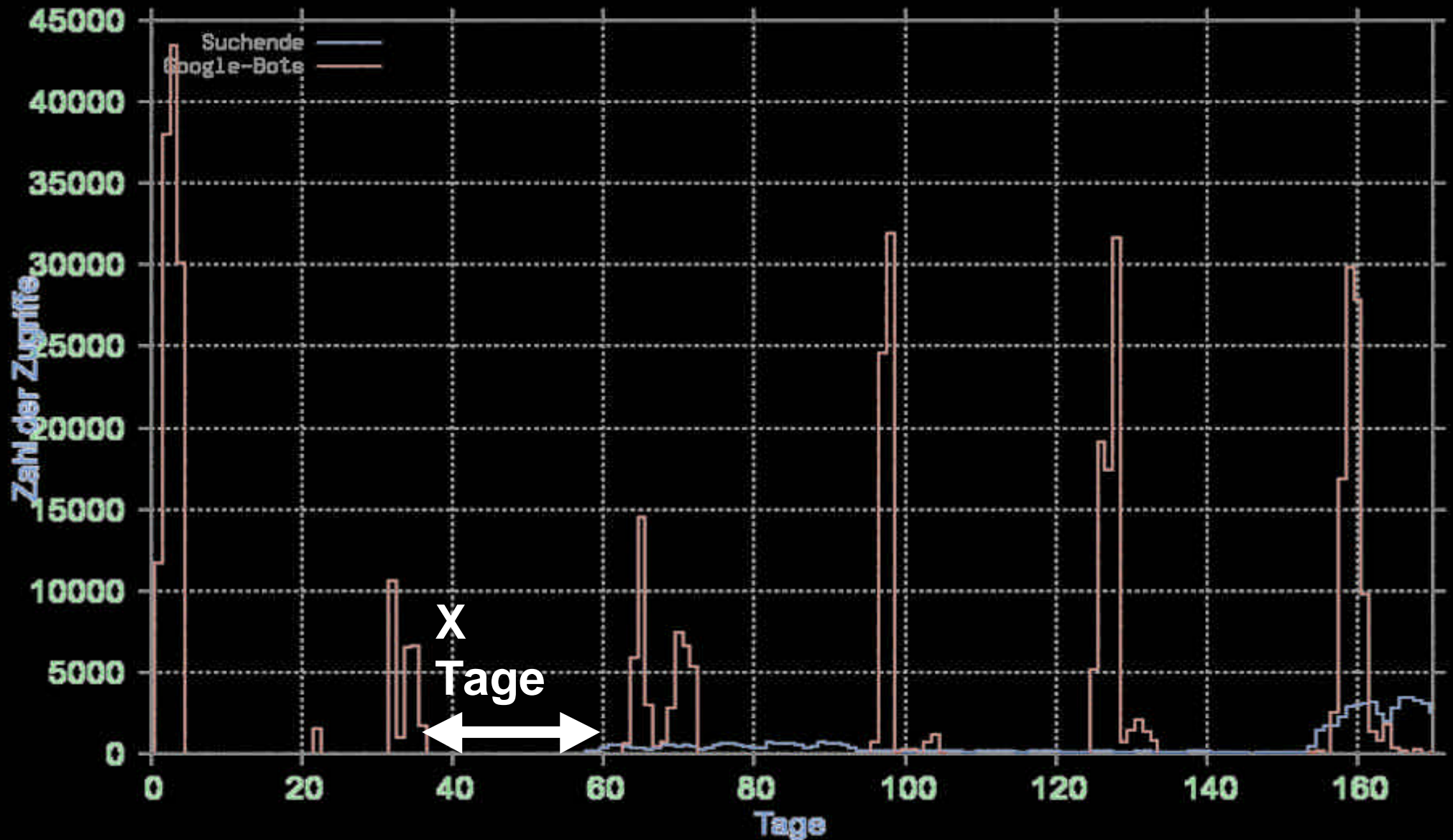
Was haben wir denn heute so vor ...

- Visualisierung ✓
- Google ✓
 - Indexierungsstrategie ✓
 - Indexierungsleistung ✓
 - Seitenzahlbegrenzung pro Server ✓
 - Dauer bis zur Findbarkeit in der Suche
 - "Tiefe" der Indexierung
 - Dateigroessen
 - Einfluss auf die Besuchshaeufigkeit der Bots
- fast, msn, neofonie
- Was lernen wir daraus

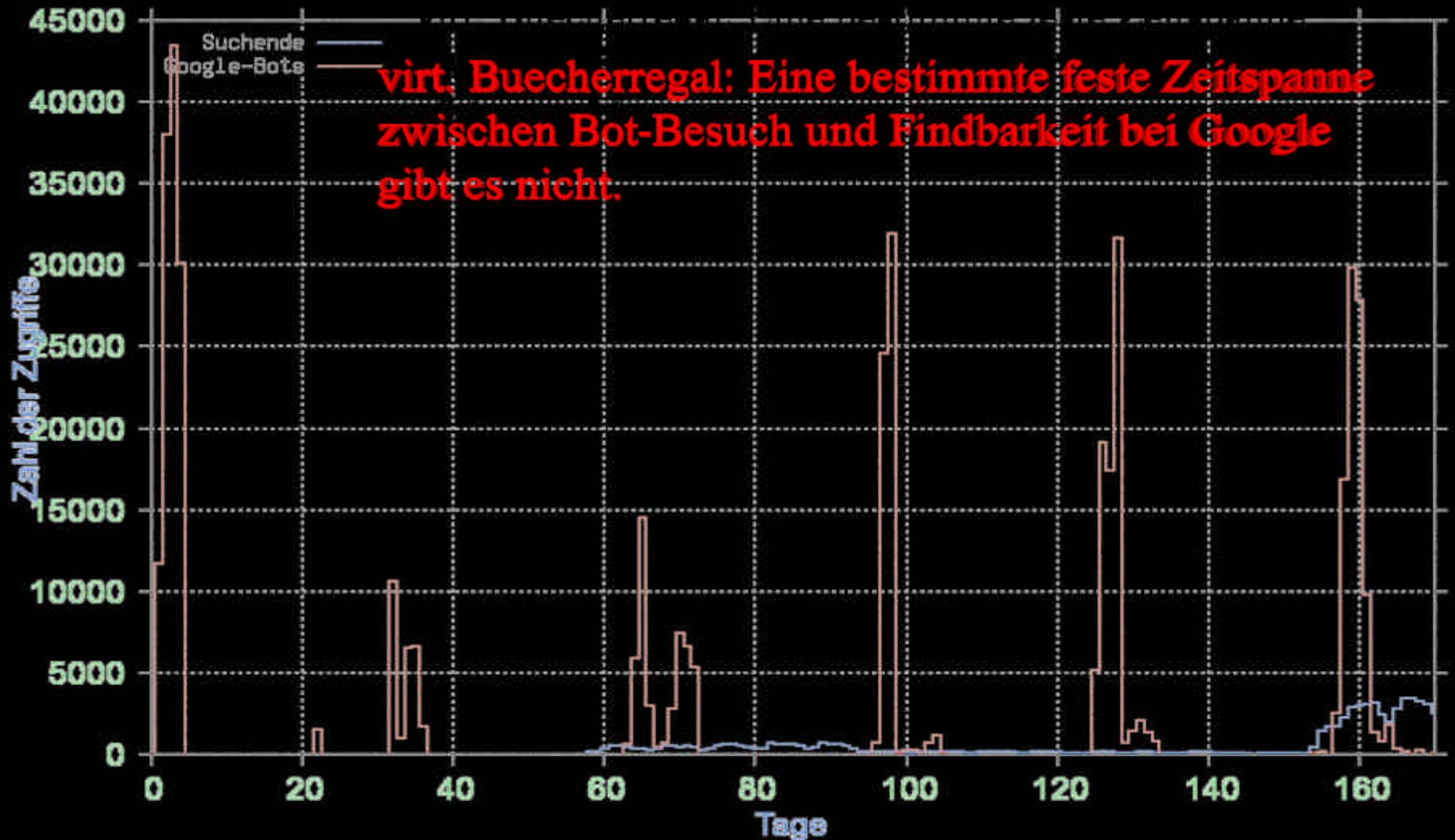
Google-Bots und Treffer durch Google-Suchende



Google-Bots und Treffer durch Google-Suchende



Google-Bots und Treffer durch Google-Suchende



"Tiefe" der Indexierung

- Geruecht:

- Google indexiert nur bis zu einer bestimmten Tiefe.
- Stimmt das aus Sicht des virt. Buecherregals?

- url:

- <http://kirke.hbz-nrw.de/>
- Stufe 1: dcb/
- Stufe 2: Alle_003/
- Stufe 3: Buecher_13/
- Stufe 4: in_NRW_13/000695592.html

"Tiefe" der Indexierung

• Geruecht:

- Google indexiert nur bis zu einer bestimmten Tiefe.
- Stimmt das das Sicht des virt. Buecherregals?

• url:

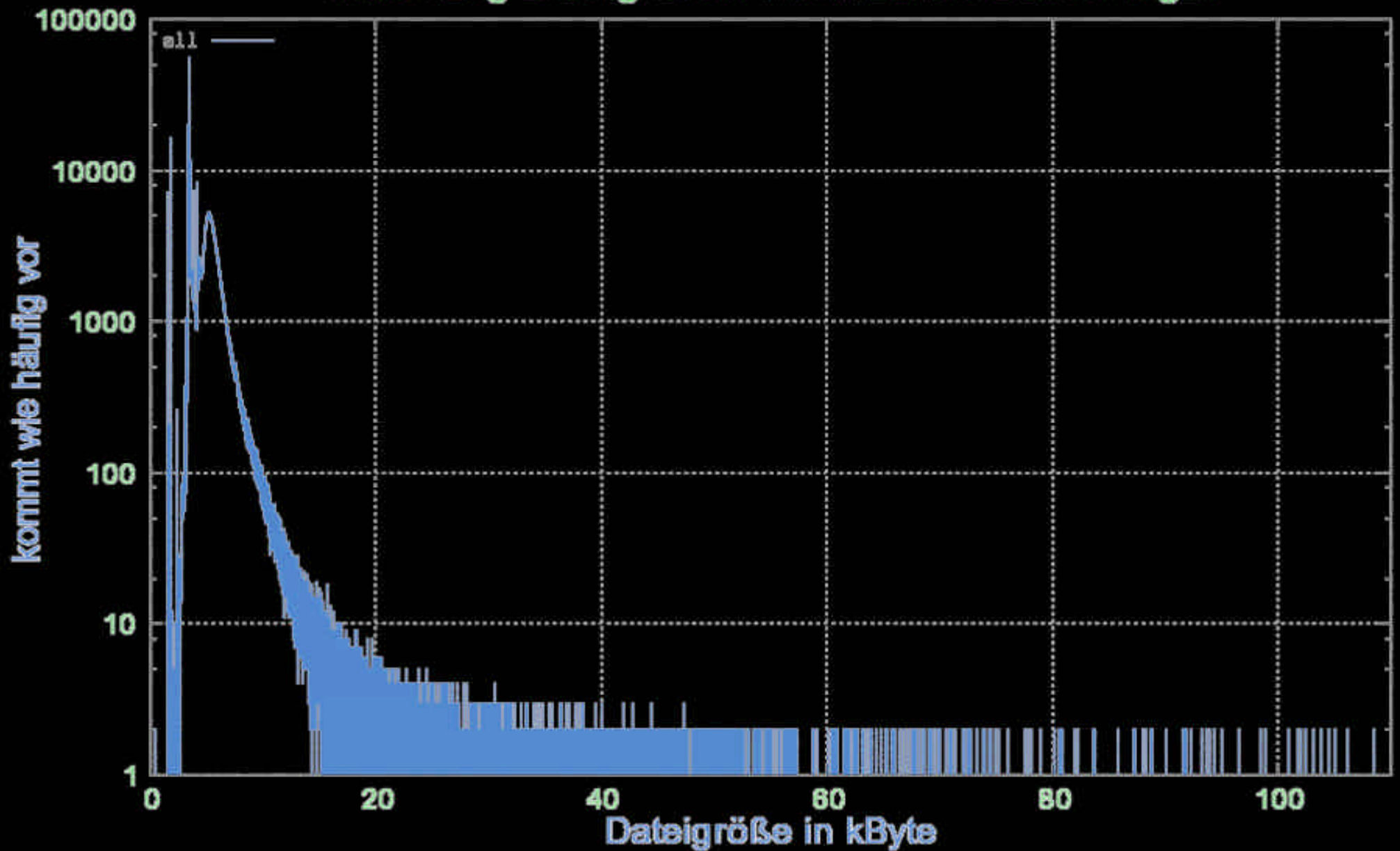
- <http://kirke.hbz-nrw.de/>
- Stufe 1: dcb/
- Stufe 2: Alle_003/
- Stufe 3: Buecher_13/
- Stufe 4: in_NRW_13/000695592.html

virt. Buecherregal: Geruecht stimmt nicht.

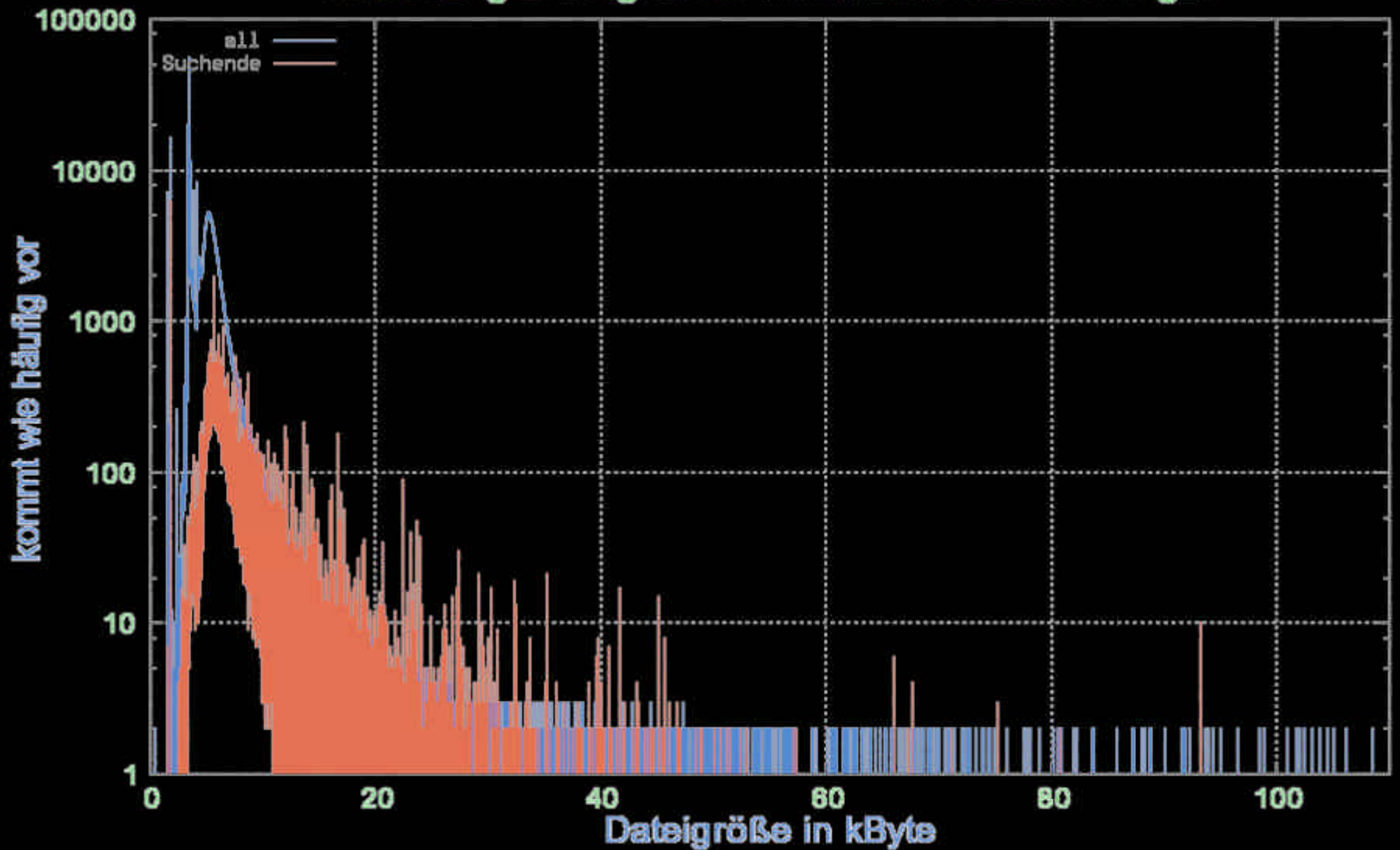
Was haben wir denn heute so vor ...

- Visualisierung ✓
- Google ✓
 - Indexierungsstrategie ✓
 - Indexierungsleistung ✓
 - Seitenzahlbegrenzung pro Server ✓
 - Dauer bis zur Findbarkeit in der Suche ✓
 - "Tiefe" der Indexierung ✓
 - Dateigroessen
 - Einfluss auf die Besuchshaeufigkeit der Bots
- fast, msn, neofonie
- Was lernen wir daraus

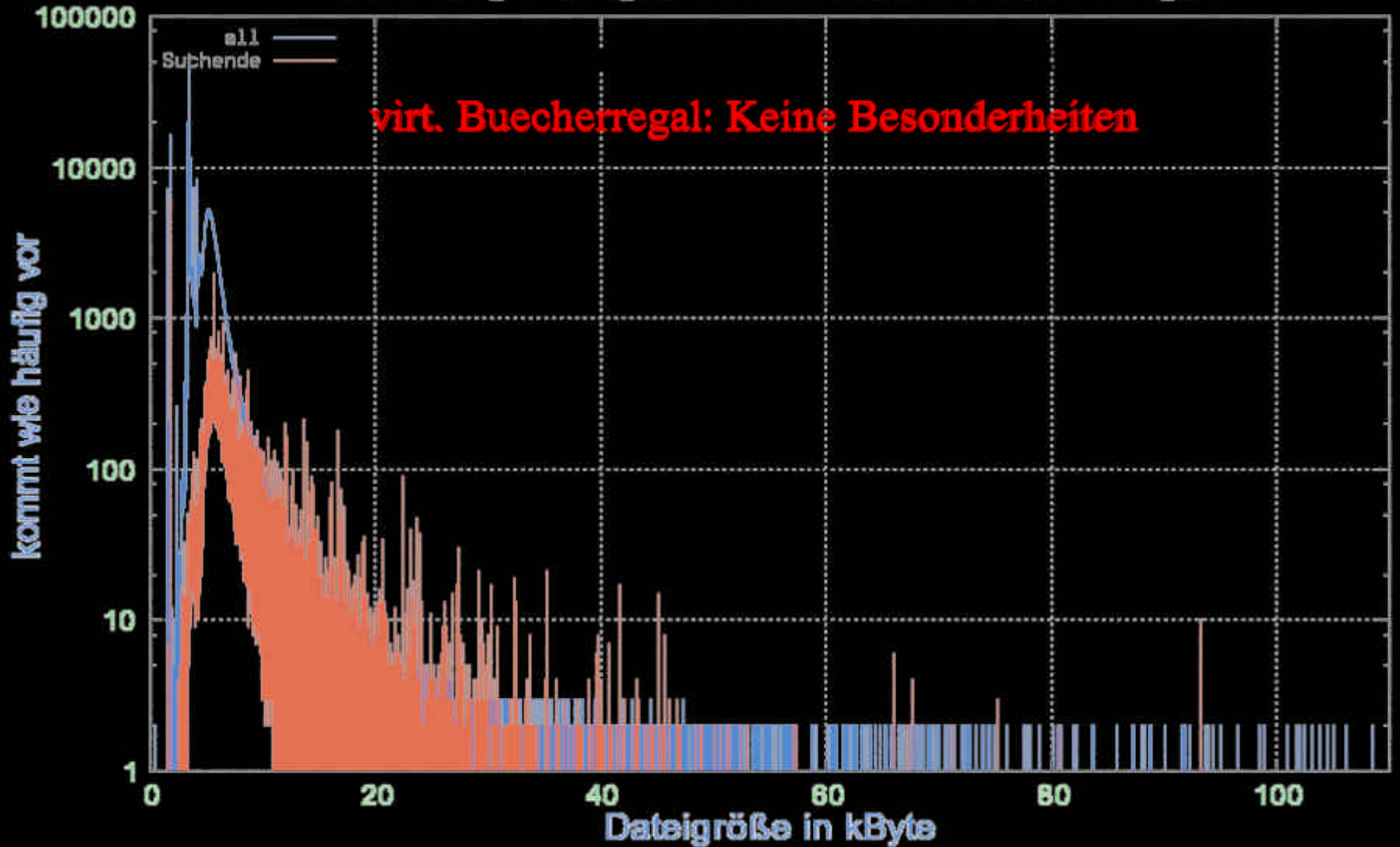
Verteilung Dateigrößen im virtuellen Bücherregal



Verteilung Dateigrößen im virtuellen Bücherregal



Verteilung Dateigrößen im virtuellen Bücherregal



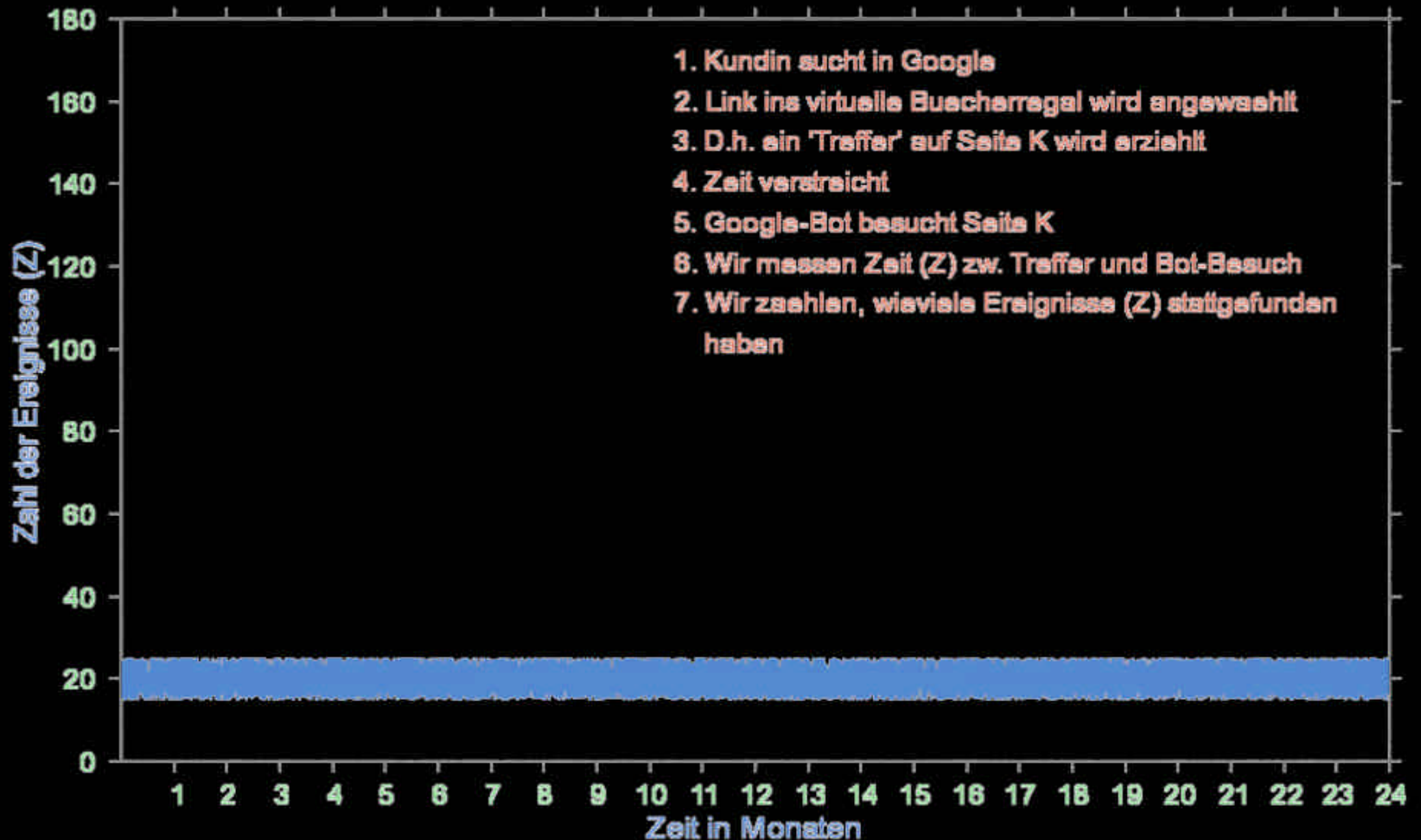
Was haben wir denn heute so vor ...

- Visualisierung ✓
- Google ✓
 - Indexierungsstrategie ✓
 - Indexierungsleistung ✓
 - Seitenzahlbegrenzung pro Server ✓
 - Dauer bis zur Findbarkeit in der Suche ✓
 - "Tiefe" der Indexierung ✓
 - Dateigrößen ✓
 - Einfluss auf die Besuchshäufigkeit der Bots
- fast, msn, neofonie
- Was lernen wir daraus

Bots:

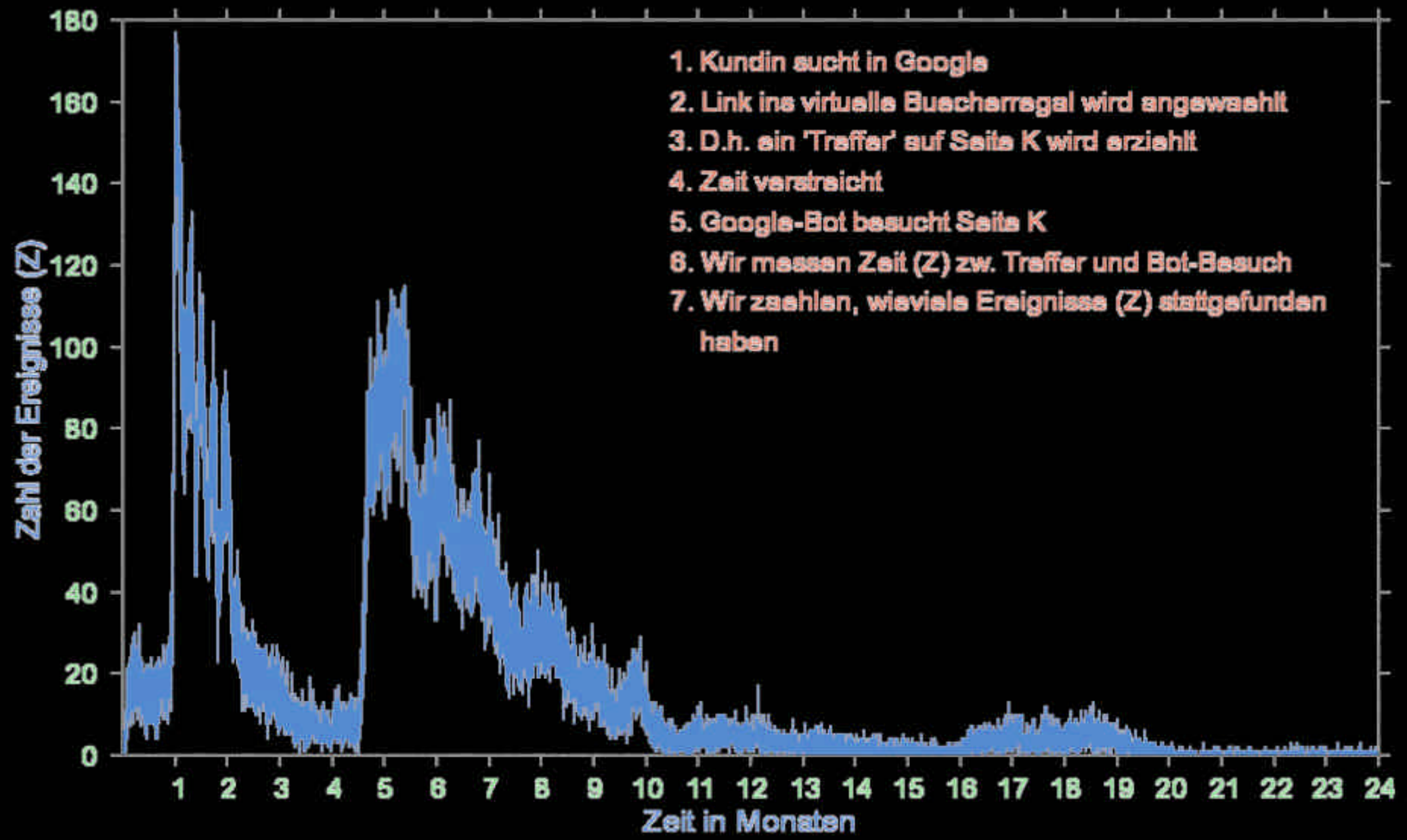
- Einfluss auf die Besuchshäufigkeit?
- Google sagt:
 - Alle 4 Mrd. Seiten werden einmal pro Monat indexiert
- virtuelles Bücherregal:
 - Wir prüfen den zeitl. Abstand zwischen einem Treffer nach einer Suche und dem folgenden Bot-Besuch.
- Was erwarten wir?

Erwartung: Zahl der Google-Bot Besuche nachdem die Seite in einer Suche gefunden wurde



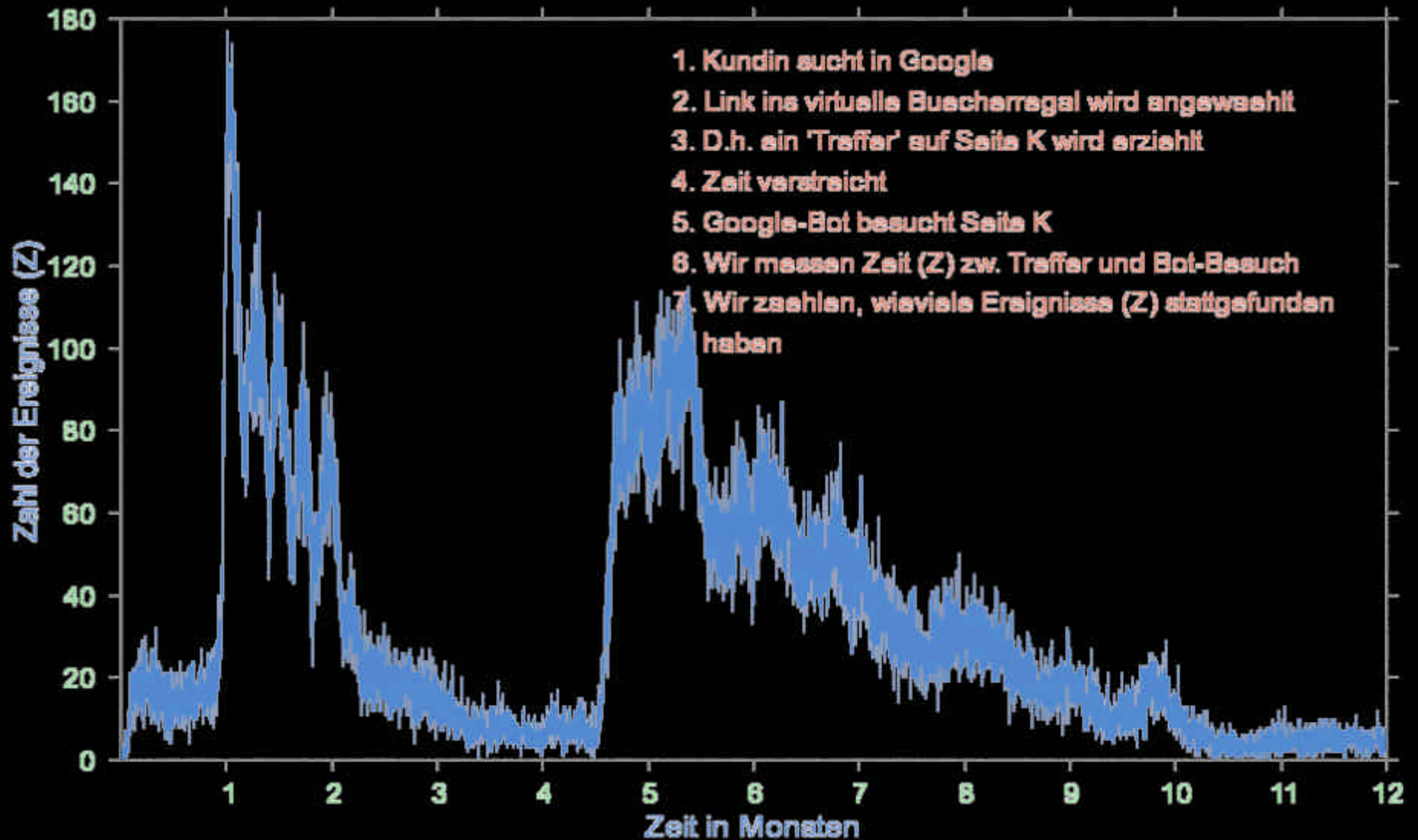
1. Kundin sucht in Google
2. Link ins virtuelle Bucherregal wird angewählt
3. D.h. ein 'Treffer' auf Seite K wird erzielt
4. Zeit verstreicht
5. Google-Bot besucht Seite K
6. Wir messen Zeit (Z) zw. Treffer und Bot-Besuch
7. Wir zählen, wieviele Ereignisse (Z) stattgefunden haben

Messung: Zahl der Google-Bot Besuche nachdem die Seite in einer Suche gefunden wurde



1. Kundin sucht in Google
2. Link ins virtuelle Bucherregal wird angewählt
3. D.h. ein 'Treffer' auf Seite K wird erzielt
4. Zeit verstreicht
5. Google-Bot besucht Seite K
6. Wir messen Zeit (Z) zw. Treffer und Bot-Besuch
7. Wir zählen, wieviele Ereignisse (Z) stattgefunden haben

Messung: Zahl der Google-Bot Besuche nachdem die Seite in einer Suche gefunden wurde



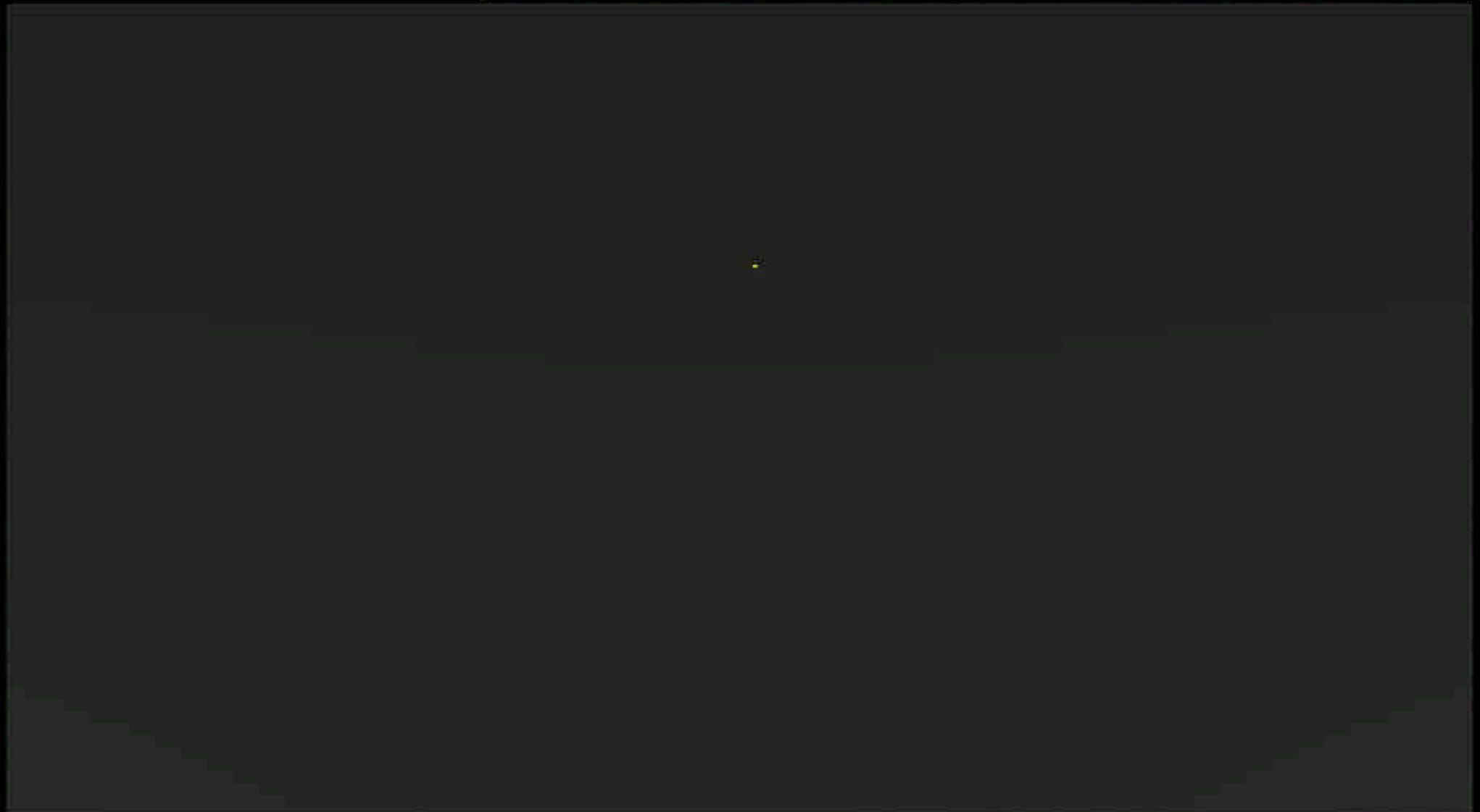
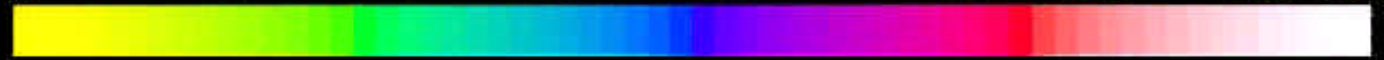
1. Kundin sucht in Google
2. Link ins virtuelle Besucherregal wird angewählt
3. D.h. ein 'Treffer' auf Seite K wird erzielt
4. Zeit verstreicht
5. Google-Bot besucht Seite K
6. Wir messen Zeit (Z) zw. Treffer und Bot-Besuch
7. Wir zählen, wieviele Ereignisse (Z) stattgefunden haben

Was haben wir denn heute so vor ...

- Visualisierung ✓
- Google ✓
 - Indexierungsstrategie ✓
 - Indexierungsleistung ✓
 - Seitenzahlbegrenzung pro Server ✓
 - Dauer bis zur Findbarkeit in der Suche ✓
 - "Tiefe" der Indexierung ✓
 - Dateigrößen ✓
 - Einfluss auf die Besuchshäufigkeit der Bots ✓
- **fast, msn, neofonie**
- **Was lernen wir daraus**

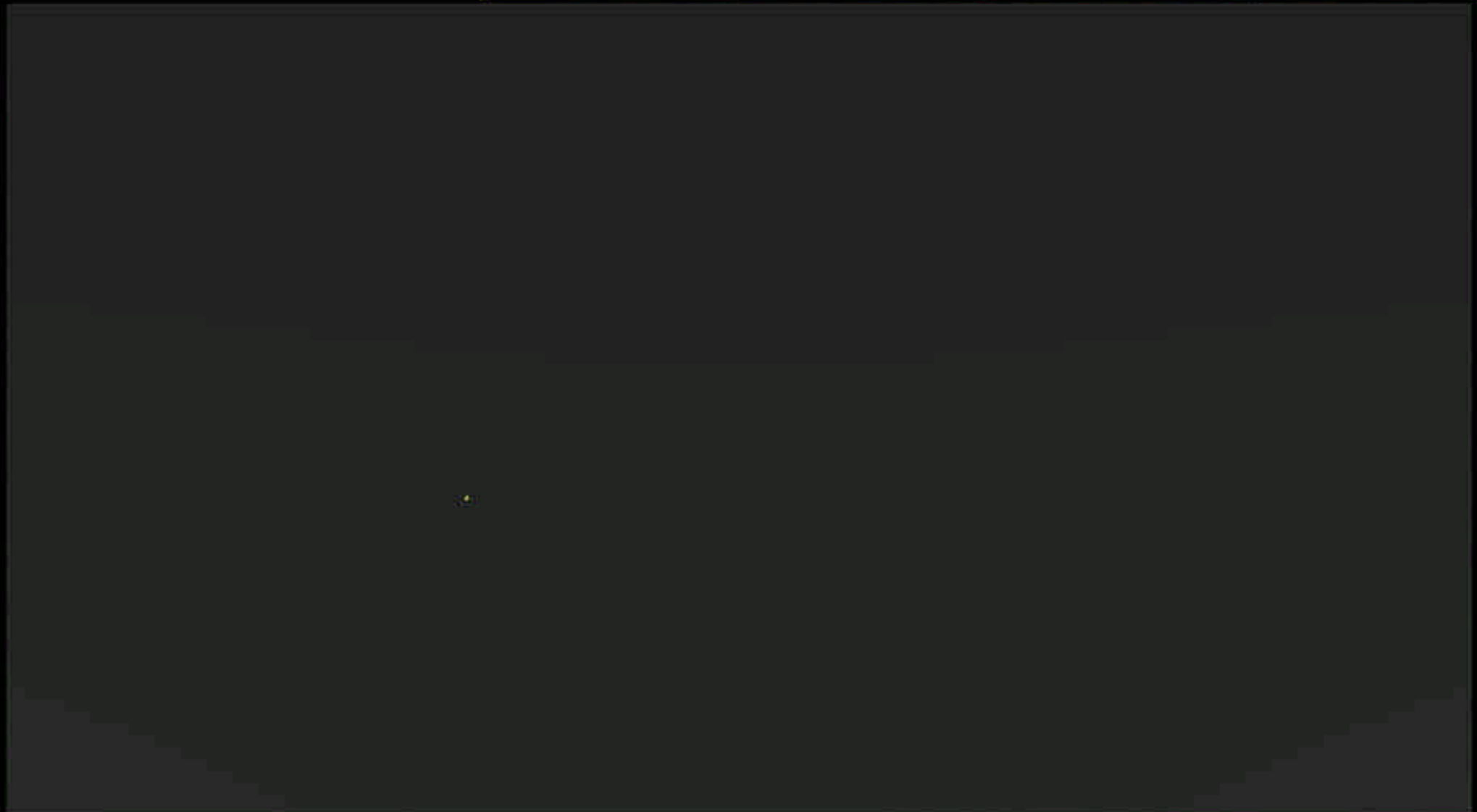
FAST-WebCrawler 16/Jul/2002

Zahl der Treffer pro Zelle: 1 2 5 10 20 60 100 200 300



msnbot 22/Jun/2003

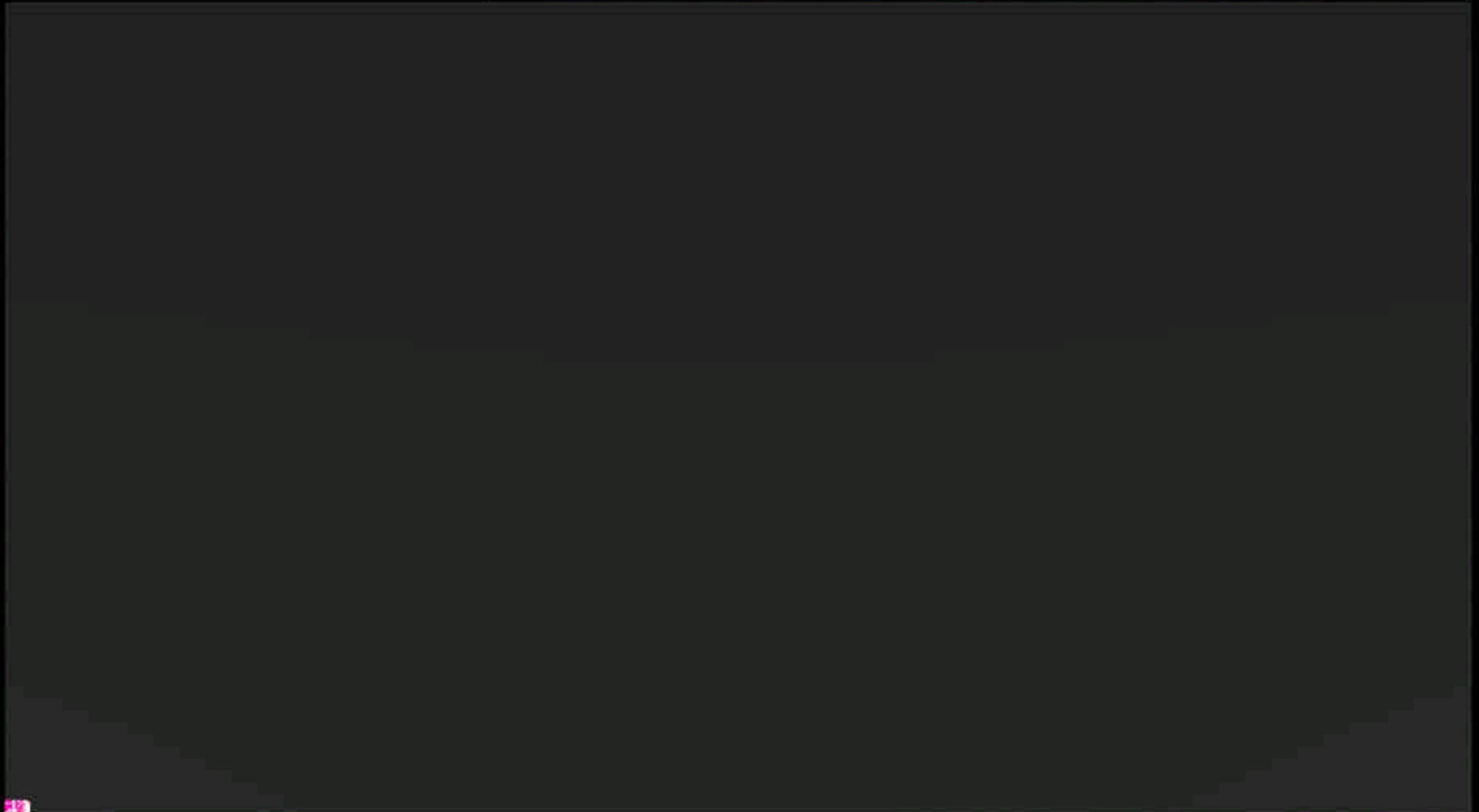
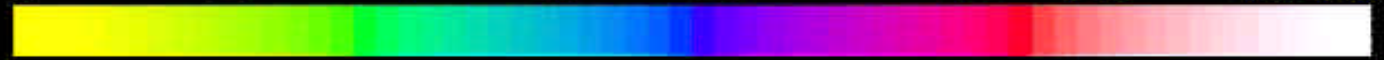
Zahl der Treffer pro Zelle: 1 2 5 10 20 60 100 200 300



Florian Seiffert, HBZ

neofonie 20/Apr/2004

Zahl der Treffer pro Zelle: 1 2 5 10 20 60 100 200 300

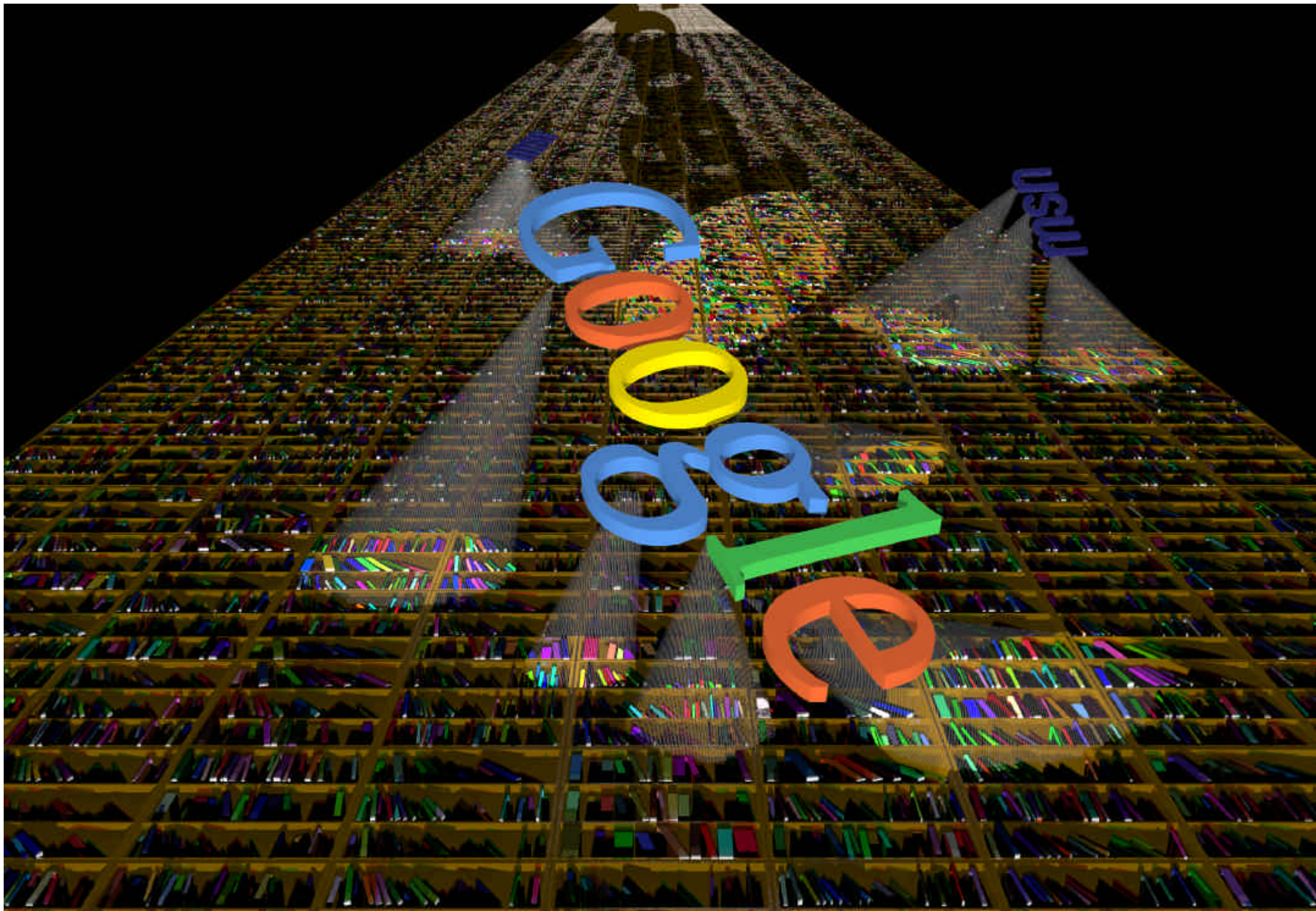


Was lernen wir daraus

- **Google**

- indexiert statistisch nicht systematisch,
- findet Seiten mit Schlagworten haeufiger,
- speichert pro Server auch mehr als 3.2 Mio Seiten,
- indexiert mindestens 4 Hierarchiestufen tief,
- aktualisiert Seiten, die gesucht werden haeufiger,

- **Suchen Sie nach Ihren Seiten !!**



Fragen?



Werkzeugkasten



- **Linux** (<http://www.suse.de>)
 - perl (<http://www.perl.com>)
 - mysql (<http://www.mysql.com/>)
 - mencoder / mplayer (<http://www.mplayerhq.hu>)
 - Persistence of Vision Raytracer (<http://www.povray.org>)
- **Windows**
 - Powerpoint
- **Vortrag unter**
 - <http://www.Florian-Seiffert.de/2004/Bonn/Inetbib2004.pdf>

"von google kommend" vom 22/Aug/2002

Zahl der Treffer pro Zelle: 1 2 5 10 20 60 100 200 300

