

Presenters: Peter Cornwell (ENS-Lyon/Data Futures) José Gonzalez (CERN) Tom Lamberty (Merve Verlag)

## FREE SEMINAR

click on link below to register



16th of October 2020 15:00-16:00 CEST



Walls Interview with Sylvère Lotringe

Lotringer: You live in East Berlin, but you are in the unusual position of being able to travel freely to the West. How much of a wall is there between East and

Müller: When I go from the Friedrichstraße chockpoint to the Zoological Garden in West-Berlin, I feel a great difference, a difference in civilization, a difference of ages, of time. There is a different time level, a different time space. You really go through a time wall. Latringer: I asked some German friends, yesterday, about East Berlin and they said: "We lave it. It's like going back to the fifties.

Müller: Most people coming from West to East Berlin keep comparing them on a horizontal level. It doesn't work that way. The problem is the misery of comparing. You just can't compare things.

Lotringer: So East Berlin is something of its own? Müller: I was very impressed by the remark of a young man who was writing an essay on my work. What was most interesting for him, he said, was this problem of another time. He remembered that he never quite unberstood why the German Wehrmacht didn't succeed in entering Moscow during the Second World War. They just stood there. They couldn't go further. He didn't believe in military or strategic reasons. He didn't be lieve in geographic reasons. He didn't believe in ideo-logical reasons. There simply was a time-wall. They ere not on the same track. This is the real problem A few years ago, I was asked by Le Monde in Paris to write something about the cultural situation in East Berlin. I tried to explain it to the French public and it wasn't easy. Then I remembered a remark Ernst Junger made. He said that you can't discuss the difference

(download master)

Last Edited: Tom Lamberty



# data **f**utures





Asset No: 10195 go save exit HTML - long HTML - paged HTML - basic

Page No: 11 / 212 go | prev | next

Status: No pending changes

er.

incoded a solves				
Walls				
Interview with Sylvère Lotringer				
Lotringer: You live in East Berlin,but you of a wall is there between East and West	are in the unusual position of being able to travel freely to the West. How much ?			
	Se checkpoint to the Zoological Garden in West-Berlin, 1 feel a great difference, ges, of time. There is a different time level, a diffe- rent time space. You really g			
Lotringer: I asked some German friends, fifties.*	yesterday,about East Berlin and they said: "Ne love it. It's like going back to the			
Müller: Most people coming from West to East Berlin keep comparing them on a horizontal level. It doesn't work that The problem is the misery of compa- ring. You just can't compare things.				
Lotringer: So East Berlin is something of its own?				
interesting for him, he said, was this prot German Wehrmacht didn't succeed in ent	ark of a young man who was writing an essay on my work. What was most blem of another time. He remembered that he neuer quite un- derstood why the tering Moscow during the Second World War. They just stood there. They couldn strategic reasons. He didn't be- lieve in geographic reasons. He didn't believe in			
Path: h3				
Final element continues next page	Blank page			
Front-matter	Footnote continues next page			
Missing/broken image	Trailing white space			
inscribed Page No: 9				
Comment:				
Please note that footland-onte contents v	will only display when the marker is high-lighted			

🔊 (\* | 👗 🖓 🖏 | B / U 🖛 | 🖩 🗃 🗐 🕊 Ω 🗄 🛨 💆 | 🚥

Heiner Müller Rotwelsch

Merve Verlag Berlin









- Peter Cornwell co-chair RDA PTTP-IG, ENS-Lyon, IMCC Westminster
- José Gonzalez Head of repository technologies, CERN
- Tom Lamberty Managing Director, Merve Verlag

email contact form: <u>www.data-futures.org/pages/contact</u>



### FRANKFURTER BUCHMES 14 - 18 OCTOBER 2020







 open reading online (albeit with access control options) using presentationcapable IIIF service supporting multiple standards-based readers

- <u>catalogue.merve.hasdai.org</u>

- print-on-demand of volumes which are no longer in print, plus custom reader print service — whether current catalog or back-list
- support for researchers via standards-based interfaces to library systems, creation, maintenance and publishing of annotations and full corpus search







- preserve the publisher: back-catalog (born-analog 1970-1990's) plus current books (born-digital 1990's-today) redelivered as agnostic HTML, multiple styled outputs as well as book metadata, using a trusted corpus-specific repository (Invenio—same tech as Zenodo)
  - preserve historic scans, PDFs, raw OCR and original digital sources including markup
  - maintain and preserve new digital resource: agnostic HTML, styled printready outputs incl. e.g. custom readers
- sustainable archive of publisher business administration documents, including agreements and licenses, correspondence with agents such as authors, editors and translators as well as distribution data like e.g. reviews







# transform existing books

- workflow-1: accession of *back-catalog* scans to produce page impressions and manage OCR data using MongoDB-based *freizo* platform
  - fix recognition errors, remove end-of-line hyphenation and make flexible linked foot/end note definitions, nested section styles and emphasis
  - generate HTML and IIIF service for reading and research applications
- workflow-2 : redelivery of digitally originated current-catalog books
  - page impressions rendered from digital sources
- both workflows produce common outputs—ASCII, un-styled HTML & page imagery





# transform existing books

- produce new virtual catalog of technology-agnostic texts
  - edit and maintain historic books
- side, together with page continuation and footnote/reference controls

ring. You just can't compare things.

- automated assistance to correct recognition errors, transform hyphenation and

- create style-able HTML references and footnotes, electronic ToCs

page impression, initial HTML and editor, plus heat-mapped OCR delivered side-by-

```
Müller: Most people coming from West to East Berlin
keep comparing them on a horizontal level. lt doesn't
work that way. The problem is the misery of compa-
Lotringer: So East Berlin is something of its own?
Müller: I was very impressed by the remark of a young
man who was writing an essay on my work. What was
most interesting for him, he said, was this problem of
another time. He remembered that he neuer quite un-
derstood why the German Wehrmacht didn't succeed in
entering Moscow during the Second World War. They
just stood there. They couldn't go further. He didn't
believe in military or strategic reasons. He didn't be-
lieve in geographic reasons. He didn't believe in ideo-
logical reasons. There simply was a time-wall. They
```





Walls Interview with Sylvere Lotringer

Lotringer: You live in East Berlin, but you are in the unusual position of being able to travel freely to the West. How much of a wall is there between East and West?

Müller: When I go from the Friedrichstraße checkpoint to the Zoological Garden in West-Berlin, I feel a great difference, a difference in civilization, a difference of ages, of time. There is a different time level, a different time space. You really go through a time wall.

Lotringer: I asked some German friends, yesterday, about East Berlin and they said: "We love it. It's like going back to the fifties."

Müller: Most people coming from West to East Berlin keep comparing them on a horizontal level. It doesn't work that way. The problem is the misery of comparing. You just can't compare things.

Lotringer: So East Berlin is something of its own?

Müller: I was very impressed by the remark of a young man who was writing an essay on my work. What was most interesting for him, he said, was this problem of another time. He remembered that he never quite understood why the German Wehrmacht didn't succeed in entering Moscow during the Second World War. They just stood there. They couldn't go further. He didn't believe in military or strategic reasons. He didn't believe in geographic reasons. He didn't believe in ideological reasons. There simply was a time-wall. They were not on the same track. This is the real problem. A few years ago, I was asked by Le Monde in Paris to write something about the cultural situation in East Berlin. I tried to explain it to the French public and it wasn't easy. Then I remembered a remark Ernst Junger made. He said that you can't discuss the difference

2) (

Headi

### Walls

Interview with Sylvère Lotringer

Lotringer: You live in East Berlin, but you are in the unusual position of being able to travel freely to the West. How much of a wall is there between East and West?

fifties."

Path: h

🔽 Fina Fro Mis

Inscrib

Comm

9

	🗈 🔀   В		ABC   📰 🗃	≣∎(	κΩ 🗄	🗄 Т 💆	HTML	
ng 3	<ul> <li>Styles</li> </ul>	•						

Müller: When I go from the Friedrichstraße checkpoint to the Zoological Garden in West-Berlin, 1 feel a great difference, a difference in civilization, a difference of ages, of time. There is a different time level, a different time space. You really go through a time wall.

Lotringer: I asked some German friends, yesterday, about East Berlin and they said: "Ne love it. It's like going back to the

Müller: Most people coming from West to East Berlin keep comparing them on a horizontal level. It doesn't work that way. The problem is the misery of compa- ring. You just can't compare things.

Lotringer: So East Berlin is something of its own?

Müller: I was very impressed by the remark of a young man who was writing an essay on my work. What was most interesting for him, he said, was this problem of another time. He remembered that he neuer quite un- derstood why the German Wehrmacht didn't succeed in entering Moscow during the Second World War. They just stood there. They couldn't go further. He didn't believe in military or strategic reasons. He didn't be-lieve in geographic reasons. He didn't believe in

h3	
al element continues next page ont-matter ssing/broken image	<ul> <li>Blank page</li> <li>Footnote continues next page</li> <li>Trailing white space</li> </ul>
ed Page No: 9	
ent:	

Please note that foot/end-note contents will only display when the marker is high-lighted





# high-availability IIIF service for online access

- IIIF book manifests generated by *freizo* with tables of contents
  - Presentation protocols ... open, standards-based APIs
- resources currently evaluated in bio-diversity community

- support multiple reader and research applications using both Image and

- 131 institutions are now members of growing IIIF consortium: existing and future IIIF-compatible applications include image viewers, and research instruments; much larger numbers of applications use IIIF image fragments

 multiple IIIF services generated automatically from conservation-quality master imagery provide parallel performance and automatic fail-over reliability for key

 hasdai-supported Merve IIIF service supports institutional guarantee to manage long-term costs through shared maintenance of multiple near-line copies





### $\blacksquare$ $\blacksquare$ $\blacksquare$ Philosophien der Literatur

**N** 

밢

### # Front matter

# Erste Vorlesung [23.04.2002]

Ø Einführung

- 1 Griechenland
- 2 Philosophie

Zweite Vorlesung [30.04.2002]

Dritte Vorlesung [07.05.
 2002]

Vierte Vorlesung [14.05.2002]

3 Mittelalter

Fünfte Vorlesung [21.05.2002]

4 Eine literaturgeschichtliche Revolution

16/10/2020



# Vorlesungen müssen

genau das sagen

JJYJJJJOUIN

# was in keinem Buch

Philosophien der Literatur

Merve

# Merve Verlag

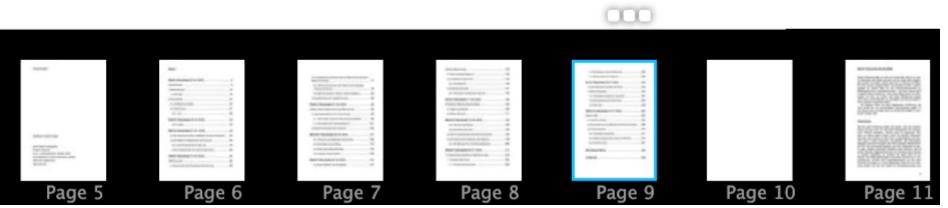
### $\blacksquare$ $\blacksquare$ Philosophien der Literatur

[07.00.2002]	Q)		27
5 Kants Nachfolger	副		7.
Achte Vorlesung [11.06.2002]			7.
6 Friedrich Wilhelm Georg Hegel		Zv	۷Ö
Neunte Vorlesung [18.06.2002]			VAC
Zehnte Vorlesung		5	3.1
[02.07.2002]	$\leq$	8	3.2
7 Literaturphilosophie als Antiphilologie		8	3.3
Elfte Vorlesung [09.07.2002]			8.
Zwölfte Vorlesung [16.07.2002]			8.
8 Nach 1945			8.
# 8.3.2 Roman Jakobson und Claude Lévi-Strauss		E	л
# 8.3.3 Jacques Lacan			
Editorische Notiz	-	Santan Manata e ana Manatan	
Literatur			
	Page 3	Page 4	

Mirador (Harvard-Stanford) IIIF viewer navigation using multi-level table of contents

16/10/2020

.3.1 DAS DASEIN IN SEINER ALLTÄGLICHKEIT	
254 A.S.2 DER URSPRUNG DES KUNSTWERKS	
259.3.3 DICHTUNG	
DLFTE VORLESUNG [16.07.2002]	
сн 1945 264	
JEAN-PAUL SARTRE	
ROLAND BARTHES ALS ÜBERGANG ZUM STRUKTURALISMUS 269	
STRUKTURALISTEN	
3.3.1 FERDINAND DE SAUSSURE	
3.3.2 ROMAN JAKOBSON UND CLAUDE LEVI-STRAUSS	
3.3.3 JACQUES LACAN	







# internationally-supported data preservation infrastructure

- search and networking facilities developed by the global community

  - comprehensive and tailorable search tools: Elasticsearch
- - Invenio in-band maintenance of book sources and metadata
  - uploading of publisher business administration documents
  - access controls for selective reader copyright protection and work-inprogress

 Invenio-3.3 repository generated by freizo provides public access and large-scale - mature services for connection with research community standards

2021 InvenioRDM release enables maintenance of catalog & archive by publisher





# Findable, Accessible, Interoperable, Reusable

- generated from book redeliveries
- significantly—IIIF manifests
- Merve records

concurrent full corpus searches via independent MongoDB engine using ASCII text

• serializer infrastructure (currently provided by freizo; built in to InvenioRDM) for generation of LoC metadata equivalents such as MODS and MARC and,

OAI Protocol for Metadata Harvesting enables external organizations to discover





# research interfaces & scholarly annotation

- linked to contributors' ORCIDs
  - page interactively
  - OADM dialects, from WADM

 sustainable reference IIIF service for current and back book catalogs, with full presentation protocol layer, supporting creation and maintenance of annotation

- enables multiple existing and future workflows to link annotations to the original

- preservation of canonical WADM annotation records using repositories such as Zenodo promotes collaboration and publishing and, significantly, serialization of multiple viewer software-specific annotation representation 'flavors', such as





### Annotations

Die Kunst, Listen zu erstellen (I)

Weiblichkeit in der Schrift (I)

Rotwelsch (12)

Geschwindigkeit und Politik (1)

Kleine Schriften (2)

Das Geschlecht, das nicht eins ist (3)

builds on existing support for WADM annotation: https://de.merve.de in 2018

2

Ę

16/10/2020

### Rotwelsch, 1982

download annotations as WADM

p11, a0 (Merve Digital Edition) **You really go through a time wall.** (vgl. Ernst Jünger: An der Zeitmauer (1959))

p13, a0 (Merve Digital Edition)

### It's an artificial freedom, an artificial free space for ideology, for the arts and for literature.

(The term 'artificial' would be worthwhile a separate intervention)

p49, a0 (Merve Digital Edition) (He was the first to see/say this)

p107, a0 (Merve Digital Edition) (Das ist mal eine andere 'Traumdeutung')





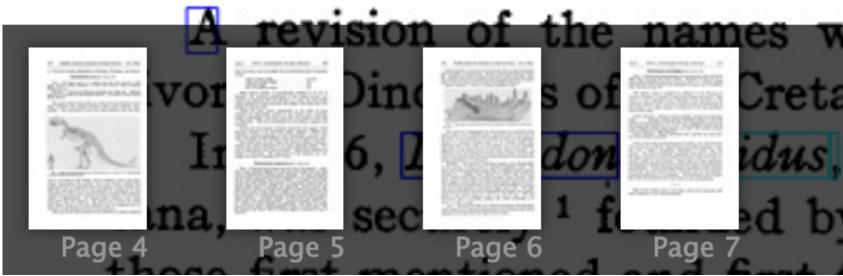
I also briefly characterize as **Dynamosaurus** another carnivorous dinosaur, with dermal plates, found by Mr. Brown in 1900. The carnivorous group has hitherto been considered as belonging to the single genus Dryptosaurus, but it is probably little less diversified than

D Plazi, a23:

taxonomicName - \_evidence < catalogs > \_step Reptilia>family< Tyrannosauridae>genus< All GBIF>kingdom< Animalia>order< Dinosauria:

D Plazi, a22:

saurus, is represented in the British Columbia skulls hitherto described as Dryptosaurus.



16/10/2020

团 Edit 面 Delete	Cera-
	by par-
o< genera>authority< Osborn, 1905>class< Ibertosaurus>higherTaxonomySource<	search. nder it
<pre>&gt;phylum&lt; Chordata&gt;rank&lt; genus&gt;</pre>	inodon
团 Edit 面 Delete	enus of
	lberto-

## I. NOMENCLATURE.

which have been applied to the Car-Cretaceous appears to be necessary. idus, from the Judith River Beds of Moned by Leidy<sup>2</sup> on Megalosaurian teeth, and





# what's planned in 2021

- single-click creation and versioning of WADM annotation collections in Zenodo, with automated DOI generation
- search API plus print-on-demand fulfillment from Invenio catalog
- Kittler Special Edition workflows producing Collected Works print and open reading volumes plus Invenio repositories of non-printable media
- Fotomuseum Winterthur Still Searching Blog Special Edition
  - print and open reading volumes plus Invenio repositories
- preservation of Merve's historical documents, including agreements and licences and correspondence with agents such as authors, editors and translators

