

Cluster of Excellence
ASIA AND EUROPE
IN A GLOBAL CONTEXT



**UNIVERSITÄT
HEIDELBERG**
ZUKUNFT
SEIT 1386

Wege zur Erschließung der frühen chinesischen Presse:

Early Chinese Periodicals Online (ECPO)

Matthias Arnold | Uni Heidelberg · HCTS · HRA

Research data – Chinese periodical press

- First decades of the 20th century
- Understudied, but dominated the contemporary print market and provide access to the "actual culture" (R. Williams, 1961)

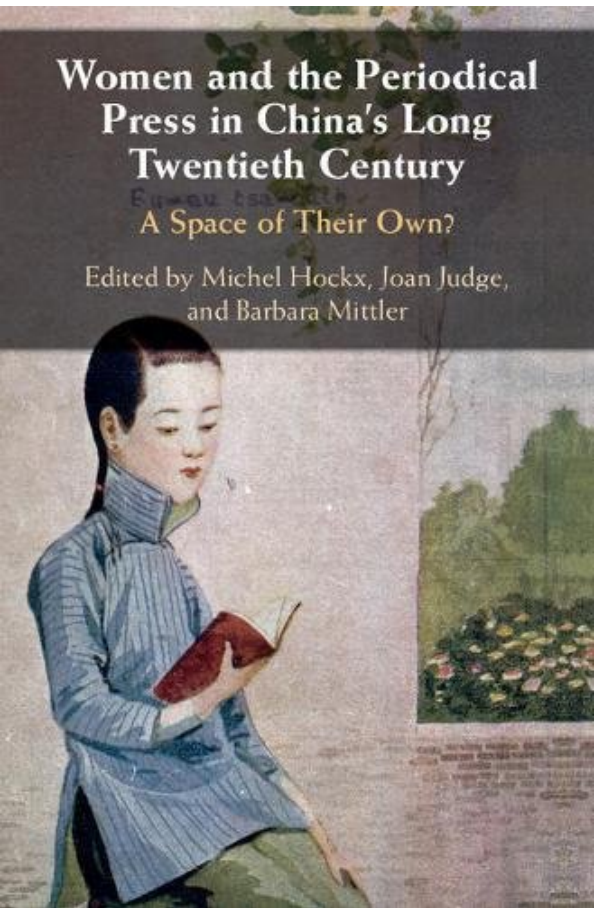
Challenges

- Physically dispersed, often poorly preserved
- Voluminous (full runs, daily, up to >30 years)
- Multi-generic and intellectually demanding

Approach

- Multi-disciplinary team, >10 researchers, 2 PhD theses
- *Women and the Periodical Press in China's Global Twentieth Century: A Space of Their Own?* Ed. by Joan Judge, Barbara Mittler and Michel Hockx, Cambridge University Press, 2018.

Database Early Chinese Periodicals Online (ECPO)



Women and the Periodical Press in China's Long Twentieth Century A Space of Their Own?

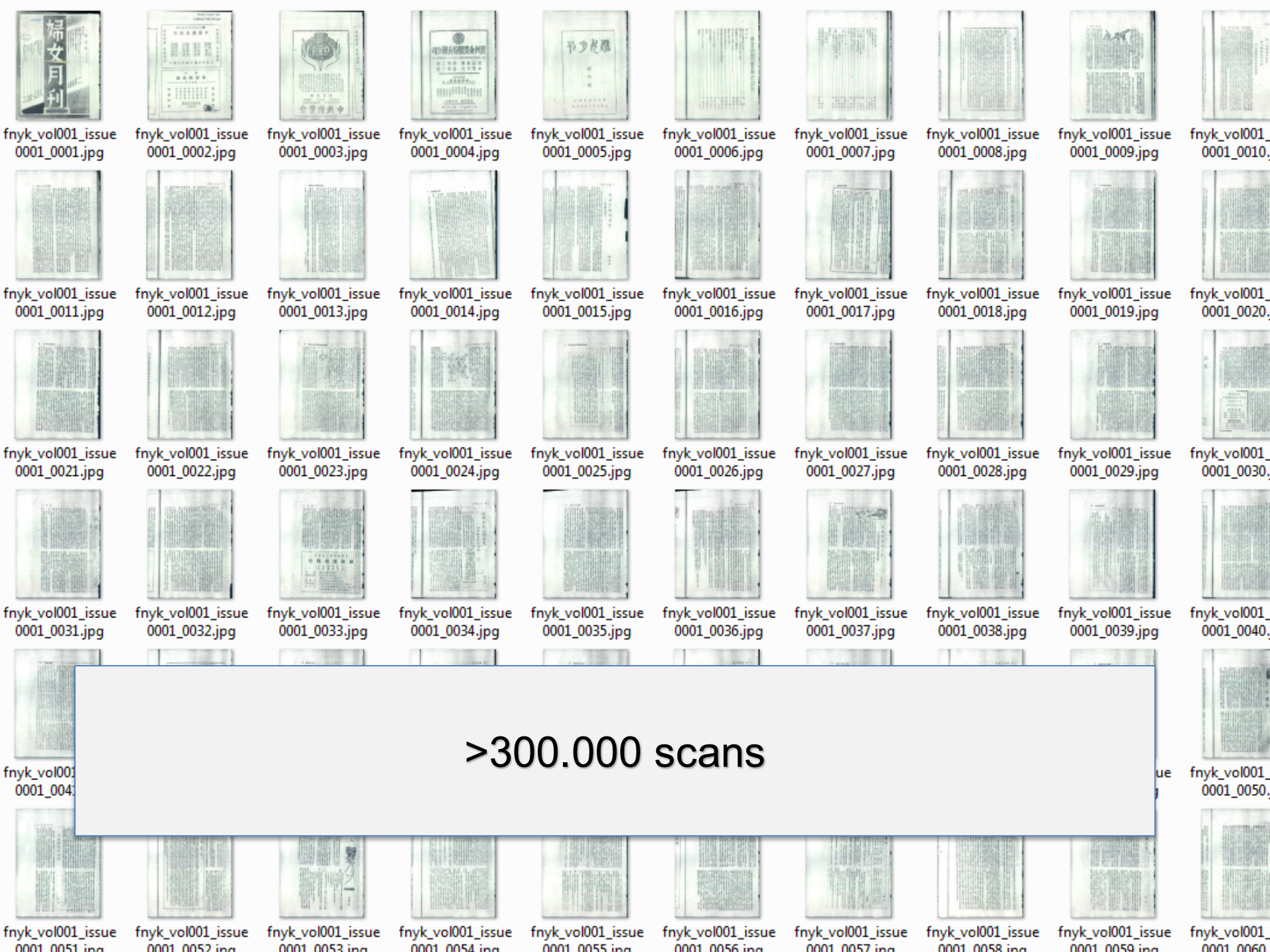
Edited by Michel Hockx, Joan Judge,
and Barbara Mittler

TYPES OF PUBLICATIONS (239)

FILM	LITERATURE	XIAOBAO	POLITICAL
GENDERED	RELIGION RELATED	FAMILY LIFE	MISCELLA- NEOUS
LIFESTYLE	PICTORIALS	MEDICINE	FASHION
YOUTH	ART	FOREIGN PRESS	

<https://uni-heidelberg.de/ecpo>





>300.000 scans

fnyk_vol001_issue_0001_0001.jpg fnyk_vol001_issue_0001_0002.jpg fnyk_vol001_issue_0001_0003.jpg fnyk_vol001_issue_0001_0004.jpg fnyk_vol001_issue_0001_0005.jpg fnyk_vol001_issue_0001_0006.jpg fnyk_vol001_issue_0001_0007.jpg fnyk_vol001_issue_0001_0008.jpg fnyk_vol001_issue_0001_0009.jpg fnyk_vol001_issue_0001_0010.jpg

fnyk_vol001_issue_0001_0011.jpg fnyk_vol001_issue_0001_0012.jpg fnyk_vol001_issue_0001_0013.jpg fnyk_vol001_issue_0001_0014.jpg fnyk_vol001_issue_0001_0015.jpg fnyk_vol001_issue_0001_0016.jpg fnyk_vol001_issue_0001_0017.jpg fnyk_vol001_issue_0001_0018.jpg fnyk_vol001_issue_0001_0019.jpg fnyk_vol001_issue_0001_0020.jpg

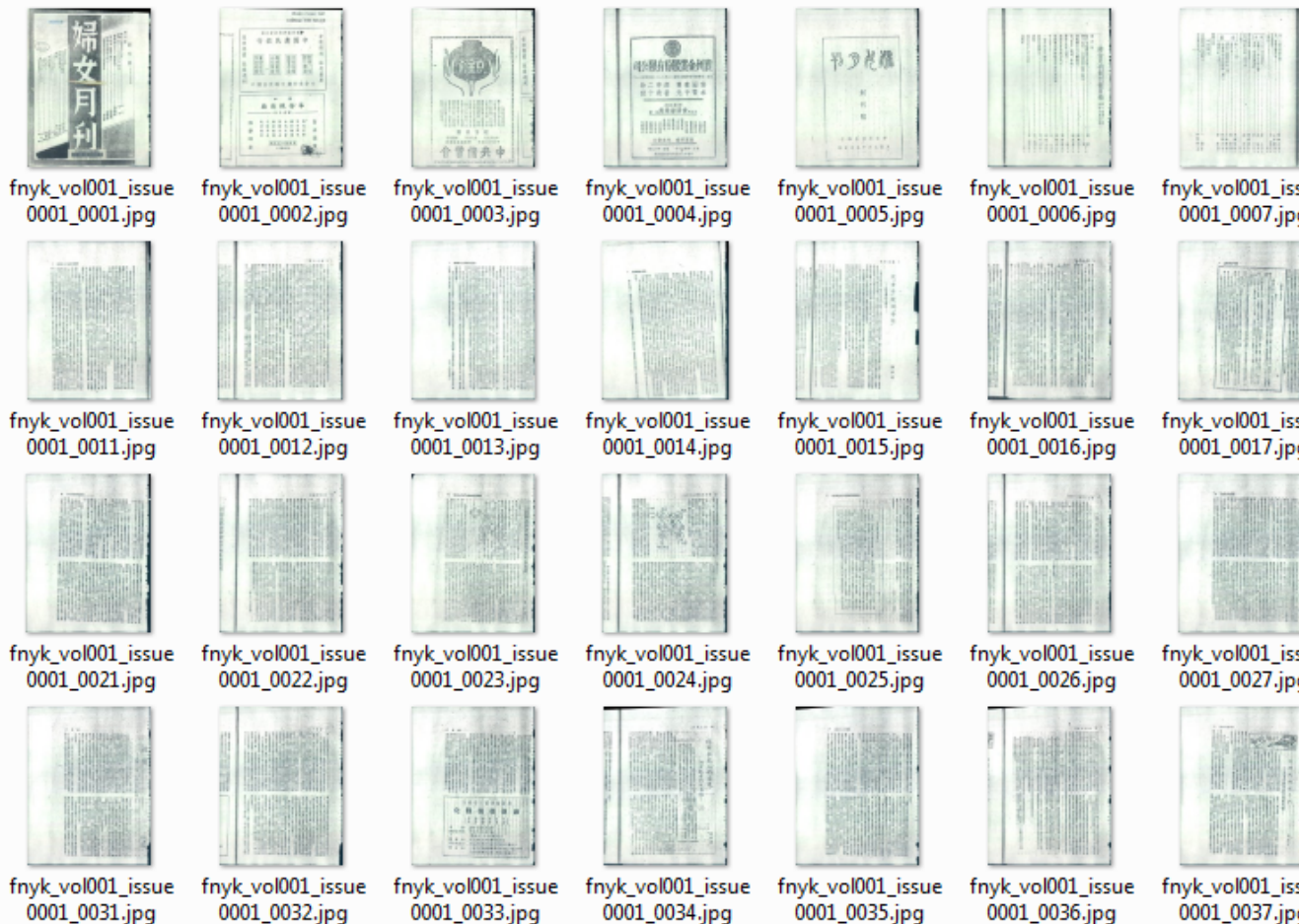
fnyk_vol001_issue_0001_0021.jpg fnyk_vol001_issue_0001_0022.jpg fnyk_vol001_issue_0001_0023.jpg fnyk_vol001_issue_0001_0024.jpg fnyk_vol001_issue_0001_0025.jpg fnyk_vol001_issue_0001_0026.jpg fnyk_vol001_issue_0001_0027.jpg fnyk_vol001_issue_0001_0028.jpg fnyk_vol001_issue_0001_0029.jpg fnyk_vol001_issue_0001_0030.jpg

fnyk_vol001_issue_0001_0031.jpg fnyk_vol001_issue_0001_0032.jpg fnyk_vol001_issue_0001_0033.jpg fnyk_vol001_issue_0001_0034.jpg fnyk_vol001_issue_0001_0035.jpg fnyk_vol001_issue_0001_0036.jpg fnyk_vol001_issue_0001_0037.jpg fnyk_vol001_issue_0001_0038.jpg fnyk_vol001_issue_0001_0039.jpg fnyk_vol001_issue_0001_0040.jpg

fnyk_vol001_issue_0001_0041.jpg fnyk_vol001_issue_0001_0042.jpg fnyk_vol001_issue_0001_0043.jpg fnyk_vol001_issue_0001_0044.jpg fnyk_vol001_issue_0001_0045.jpg fnyk_vol001_issue_0001_0046.jpg fnyk_vol001_issue_0001_0047.jpg fnyk_vol001_issue_0001_0048.jpg fnyk_vol001_issue_0001_0049.jpg fnyk_vol001_issue_0001_0050.jpg

fnyk_vol001_issue_0001_0051.jpg fnyk_vol001_issue_0001_0052.jpg fnyk_vol001_issue_0001_0053.jpg fnyk_vol001_issue_0001_0054.jpg fnyk_vol001_issue_0001_0055.jpg fnyk_vol001_issue_0001_0056.jpg fnyk_vol001_issue_0001_0057.jpg fnyk_vol001_issue_0001_0058.jpg fnyk_vol001_issue_0001_0059.jpg fnyk_vol001_issue_0001_0060.jpg

- ▶ Funu shenghuo yuekan
- ▶ Funu yuekan
 - ▶ data transfer
 - ▶ images_renamed
 - ▶ volume001
 - issue001
 - issue002
 - issue003
 - issue004
 - issue005
 - issue006
 - ▶ volume002
 - issue001
 - issue002
 - issue003
 - issue004
 - issue005
 - issue006
 - ▶ volume003
 - issue001
 - issue002
 - issue003
 - issue004
 - issue005
 - issue006



40.936 issues: 46.931 articles, 20.532 images, 18.639 ads

- issue003
- issue004
- issue005
- ▶ volume006
 - issue001



長篇小說 新上海春秋 (三三五) - CHANG PIAN XIAO SHUO XIN SHANGHAI CHUN QIU (SAN SAN WU)

Article. 晶報 Jing bao, Volume 1, Issue 3843, Friday, 1939-04-21, Page 2

TITLE:	長篇小說 新上海春秋 (三三五) Chang pian xiao shuo xin Shanghai chun qiu (san san wu) Full-length novel: New Shanghai Year (335)
PUBLICATION:	晶報 Jing bao "The Crystal"
PUBLICATION TYPE:	NEWSPAPER
DATE:	21 April 1939
PAGE:	2
PAGE NUMBER AS PRINTED:	2
SPECIAL POSITION:	
SEQUENCE:	27
DOCUMENT TYPE:	Article
KEYWORDS:	law, upper class, bandit, lawyer
AGENTS:	周局長 Zhou ju zhang (mentioned in article), 唐少魯 Tang Shaolu (mentioned in article), 張延年 Zhang Yannian (mentioned in article), 楊賓秋 Yang Binqiu (mentioned in article), 沈家小姐 Shen jia xiao jie (mentioned in article), 董綬臣 Dong Shouchen (mentioned in article)
ITEM LAST MODIFIED:	
DOCUMENT URL:	http://ecpo.uni-hd.de/publications.php?magazin_id=1&isid=20&ispage=2&itemid=299&itype=2



電聲日報
DIAN SHENG RI BAO
"RADIO MOVIE DAILY NEWS"



PUBLISHING INFORMATION



- Founding Date:** 1932-05-01
- End of publication:** 1933-12-28 (uncertain)
- Frequency:** Daily.
- Format/Size:** Size of a sheet: 40 x 28 cm.
Format: 4 pages
..... Details: medium-sized quadratic pages, 2 pages per 1 side of a sheet
..... folding: vertically
- Print/Type:** - Original paper: issue no. 33 (1932-06-02): page 1: (本報) "是電影與無線電鼻祖；是彩紙四開報之創始者..."
- Scans: are generally of bad quality, i.e. the images are very dark, rendering some texts and almost all of the pictures difficult or impossible to read.
- Price:** Standard price
- issue no. 3 (1932-05-03): page 2: special (introductory?) price: 本埠零售：現特價僅售銅元六枚
- issue no. 6 (1932-05-06): page 4: standard price: 本報價目：每日四頁四分；每月大洋七角；每三個月二元（寄費在內）

Subscription price:
- issue no. 4 (1932-05-04): page 2: subscription price: 優待訂閱：直接定戶每月連寄費原價一元一角現以創刊紀念特價每月僅收七角即小洋八角
- issue no. 124 (1932-09-01): page 1: 本報自本日起、對於新訂閱者概照特價收費、每月僅七角、。。。
- Prominent agents:** 婦女日報 fu nu ri bao ()
- Summary of Content:** Columns/ sections:
1. Dian ying 電影 Movie Critic
... issue no. 0001 (1932-05-01) - last issue

Search for Agent:

exact search begin with Filter Lang

[Agents](#) | [Find Duplicates](#) | [Reset Sort Order...](#) | [Add new Agent](#) | [Export Agents as csv](#)

- 47245 results: [1-100](#), [101-200](#), [201-300](#), [301-400](#), [401-500](#), [501-600](#), [601-700](#), [701-800](#), [801-900](#), [901-1000](#), [1001-1100](#), [1101-1200](#), [1201-1300](#), [1301-1400](#), [1401-1500](#), [1501-1600](#), [1601-1700](#), [1701-1800](#), [1801-1900](#), [1901-2000](#), [2001-2100](#), [2101-2200](#), [2201-2300](#), [2301-2400](#), [2401-2500](#), [2501-2600](#), [2601-2700](#), [2701-2800](#), [2801-2900](#), [2901-3000](#), [3001-3100](#), [3101-3200](#), [3201-3300](#), [3301-3400](#), [3401-3500](#), [3501-3600](#), [3601-3700](#), [3701-3800](#), [3801-3900](#), [3901-4000](#), [4001-4100](#), [4101-4200](#), [4201-4300](#), [4301-4400](#), [4401-4500](#), [4501-4600](#), [4601-4700](#), [4701-4800](#), [4801-4900](#), [4901-5000](#), [5001-5100](#), [5101-5200](#), [5201-5300](#), [5301-5400](#), [5401-5500](#), [5501-5600](#), [5601-5700](#), [5701-5800](#), [5801-5900](#), [5901-6000](#), [6001-6100](#), [6101-6200](#), [6201-6300](#), [6301-6400](#), [6401-6500](#), [6501-6600](#), [6601-6700](#), [6701-6800](#), [6801-6900](#), [6901-7000](#), [7001-7100](#), [7101-7200](#), [7201-7300](#), [7301-7400](#), [7401-7500](#), [7501-7600](#), [7601-7700](#), [7701-7800](#), [7801-7900](#), [7901-8000](#), [8001-8100](#), [8101-8200](#)

number of results:
[10](#) [20](#) [50](#) [100](#) [500](#)

Name	Pinyin ↑	Lang	Type	Gender	edit	ECPO	FZ	Names	P-ID /N-ID	P-ID to merge into
Cooper, Jackie		English	Given Name	male	edit detail	6	1	5	8462/9632	<input type="text"/>
賈克哥根	Jiake Gegen	Chinese	Other Name, Variants						8462/8786	merge
茄基苛伯	Qieji Kebo	Chinese	Other Name, Variants						8462/42862	
賈克柯根	Jiake Kegen	Chinese	Other Name, Variants						8462/51938	
賈克珂潑	Jiake Kebo	Chinese	Other Name, Variants						8462/51939	
Frederick, Pauline		English	Given Name	female	edit detail	1	0	2	9147/9634	<input type="text"/>
保林佛特立	Baolin Foteli	Chinese	Given Name						9147/9633	merge
Gilbert, John		English	Given Name	male	edit detail	3	0	2	8983/9637	<input type="text"/>
約翰吉										merge
McL...										<input type="text"/>
萊波爾										merge
McLaglen, Victor		English	Given Name	male	edit detail	1	0	2	9159/9650	<input type="text"/>
維多麥克拉倫	Weiduo Maikelalun	Chinese	Given Name						9159/9647	merge
Negri, Pola		English	Given Name	female	edit detail	1	0	2	9158/9651	<input type="text"/>

47.245 agents, 163.408 occurrences, 15 'languages'

Opening the data silo

From static export to **dynamic data service**

- Output data using the Metadata Object Description Schema (MODS) - Open Access: <http://ecpo.uni-hd.de/api/mods/>

From static pre-rendered files to **dynamic image service**

- Implementation of International Image Interoperability Framework (IIIF) Image API <http://iiif.io/technical-details/>

From separate names to **cross-db agents service**

- Identify agent, assign names, link to authorities, structure information, feed data back to authority files (GND)



Wege zum Volltext:

Vorarbeiten

JADH 2018



定價報格

每三日出一張增刊無定
 中國境內全年二元半年
 元一角每月二角日本同
 餘外國各埠加倍均作大
 計算報費先惠郵費在內
 票加一

今日一張售大洋二分

第七百零三號

「特等」登於新聞之中價面議
 「頭等」封面分計三英寸寬
 「二英寸高為一格每格每格洋

新造 洋房 廉價

茲有坐落法租界唐家灣
 靈樞萬路水福里內新造
 高大三層樓洋房一幢汽
 車間洋式白磁浴缸面盆
 冷熱水龍頭西式活動廁
 所電燈自流井水一應俱
 全外附華式大花廳連樓
 一幢油漆精雅地位合宜
 深合公館住宅之用租金

清儀閣所藏 古器文物

全十册 夾紙一約 預十元 索樣請附 郵票六分

上海商務印書館啟

包天笑先生著

上海春秋

第二集出版

登場人物愈多
 情節愈加熱鬧
 文字愈加緊湊
 欲知海上種種奇事豔事
 不可不讀此書

每集二册 每集一元二角

上海大東書局啟

上海春秋第二集出版 - SHANGHAI CHUNQIU DI ER JI CHU BAN
Advertisement. 晶報 Jing bao, Volume 1, Issue 703, Sunday, 1925-01-04, Page 1

TITLE: 上海春秋第二集出版
Shanghai Chunqiu di er ji chu ban

PUBLICATION: 晶報 Jing bao "The Crystal"

PUBLICATION TYPE: NEWSPAPER

DATE: 04 January 1925

PAGE: 1

PAGE NUMBER AS PRINTED: 1

SPECIAL POSITION:

SEQUENCE: 1

DOCUMENT TYPE: Advertisement

PRODUCT BRAND:

ADVERTISEMENT CATEGORY: 16

KEYWORDS: books, Shanghai

AGENTS: 包天笑 Bao Tianxiao (mentioned in advertisement), 大東書局 Dadong shu ju (mentioned in advertisement)

NOTES ON ITEM: 包天笑先生著上海春秋第二集出版登場人物愈多情節愈加熱鬧文字愈加緊湊欲知海上種種奇事豔事不可不獨此書每集二冊，每集一元二角上海大東書局啟

ITEM LAST MODIFIED:

DOCUMENT URL: http://ecpo.uni-hd.de/publications.php?magazin_id=1&isid=282&ispage=1&itemid=4754&itype=4



Wege zum Volltext:

Vorarbeiten - OCR

JADH 2018

[Home](#)[Product Tour](#)[Plans and Pricing](#)[Support](#)[SLA](#)[Security](#)[Demo](#)[For Students](#)[About us](#)

Load source file (*.png or *.jpg):

jb_0016_1919-04-18_0002+0003.jpg

I agree with the [Terms of use](#)Or paste url to source file (*.png or *.jpg):

Recognition languages

Chinese Traditional

Profile

DocumentConversion

Image source

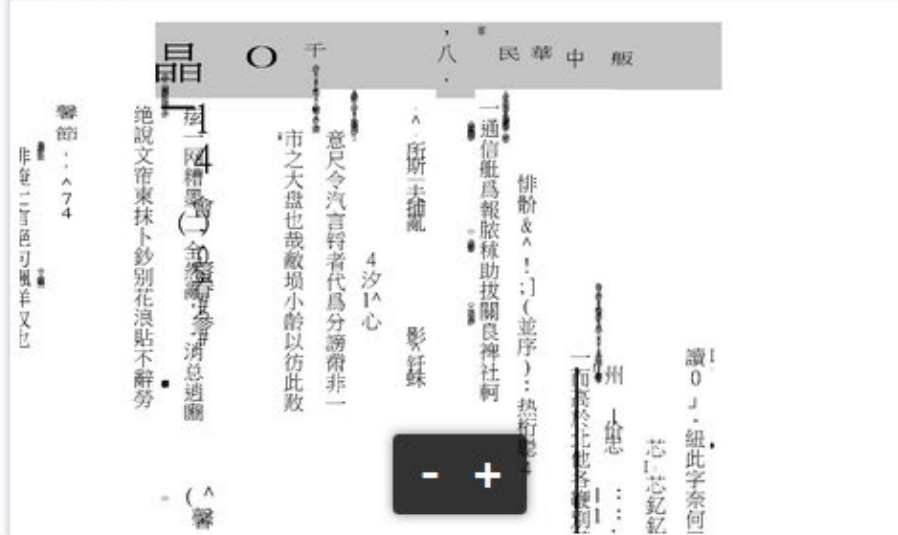
Auto

[More languages](#)

Source



Result



- Home
- Product Tour
- Plans and Pricing
- Support
- SLA
- Security
- Demo
- For Students
- About us

Load source file (*.png or *.jpg):

Browse

Or paste url to source file (*.png or *.jpg):

Recognize whole image

Recognition languages

Chinese Traditional

Profile

DocumentConversion

[More languages](#)

世哉戲填小詞以懲此輩
責者代為分謗甯非一
惡夫捕風捉影附會無
私濫屣曷翹讚之不暇

of use

Source

還捏紛紛電報
塗說文章東抹
建一團糟墨白
擲壁最能虛造
市之大意也
意且令負責
根之事張
焉所用惡惟
惡夫捕風捉影
附會無
私濫屣曷翹
讚之不暇
非小是
意遠聽
不辭勞
良裨社會
惡消聽也

排勝並行號總也
通信社為報旅藉助披關良裨社汽
符^^察一溢眾砧竝激之不暇
影附等無

Report bad result

俳體西江月（並序）惡消聽也

通信社爲報館輔助機關良裨社會
苟能偵察隱私溫犀曷翅讚之不暇
焉所用惡惟惡夫捕風捉影附會無
根之事張皇不經之言名惟徒亂人
意且令負責者代爲分謗甯非一
市之大蠹也哉戲填小詞以懲此輩
云爾

嚮壁最能虛造閉門儘可與謠非非是
是一團糟墨白全然亂了 消息道聽
塗說文章東抹 鈔郵花浪貼不辭勞
還捏紛紛電報

俳
通·信·社·
苟·能·偵·
焉·所·用·
根·之·事·
意·且·令·
市·之·大·
云·爾·
嚮·壁·最·能·
是·一·團·
塗·說·文·
還·捏·紛·

還·捏·紛·紛·電·報·
塗·說·文·章·東·抹·
是·一·團·糟·墨·白·全·然·亂·了·
嚮·壁·最·能·虛·造·閉·門·儘·可·與·謠·非·非·是·
消·息·道·聽·

俳辭一〓〓（並序）惡術聽也

通信社為報腋秬助餞關良裨社鈔
苟〓〓察一〓泓擘〓〓郊謗之不暇
忍所用惡惟惡夫捕風從影附會無
根之事張皇不經之言名佞徒亂人
意〓〓令負責寶者代為分謗甯#

市之大茲也哉敝垣小試以彷彿此依

云蜩

箔壁最能虛逝阳門儘可典謠非弗是
是一网糴墨打全然亂多消息遺艸
淦說文帘束抹。鈔郵花浪貼不辭勞

通捏紛紛訊軋

嚮壁最能虛造閉門儘可與謠非非是

是一團糟墨白全然亂了 消息道聽

塗說文章束抹 鈔郵花浪貼不辭勞

還捏紛紛電報

嚮壁最能虛造閉門儘可與謠非非是

是一團糟墨白全然亂了 消息道聽

塗說文章束抹 鈔郵花浪貼不辭勞

還捏紛紛電報

俳

通信社

苟能信

焉所用

根之事

意見令

市之大

云爾

嚮壁最能

是一團糟

塗說文章

還捏紛紛

俳辭一々（並序）惡術聽也

通信社為報腋秭助餞關良裨社鈺

苟々察Ⅰ泓庵3. 郊謗之不暇

忍所用惡惟惡夫捕風從影附會無

根之事張皇不經之言名佞徒亂人

意。令負責寶者代為分謗甯# |

市之大茲也哉敝垣小試以彷彿依

云蜩

箔壁最能虛逝阳門儘可典謠非弗是

是一网糴墨打全然亂多消息道艸

淦說文帘束抹。鈔郵花浪貼不辭勞

通捏紛紛訊軋

嚮壁最能虛造閉門儘可與謠非非是

是一團糟墨白全然亂了 消息道聽

塗說文章束抹 鈔郵花浪貼不辭勞

還捏紛紛電報

嚮壁最能虛造閉門儘可與謠非非是

是一團糟墨白全然亂了 消息道聽

塗說文章束抹 鈔郵花浪貼不辭勞

還捏紛紛電報

俳

俳

通信社

苟能信

焉所用

根之事

意見令

市之大

云爾

通信社

苟能信

焉所用

根之事

意見令

市之大

云爾

嚮壁最能虛造閉門儘可與謠非非是

是一團糟墨白全然亂了 消息道聽

塗說文章束抹 鈔郵花浪貼不辭勞

還捏紛紛電報

ca. 63% correctly recognized

本報具有廿一年的歷史是小型報紙其祖籍路普通全國刊登廣告最有力

法界懸旗糾紛

今日可復業

被捕市民昨日均已釋放

法界各界昨日已復業，法界各界昨日已復業，法界各界昨日已復業...

河內 航線已通，河內航線已通，河內航線已通...

外交界動靜，外交界動靜，外交界動靜...

本市新聞，本市新聞，本市新聞...

本市新聞，本市新聞，本市新聞...

本市新聞，本市新聞，本市新聞...

本市新聞，本市新聞，本市新聞...

本市新聞，本市新聞，本市新聞...

本市新聞，本市新聞，本市新聞...

本市新聞，本市新聞，本市新聞...

本市新聞，本市新聞，本市新聞...

本市新聞，本市新聞，本市新聞...

本市新聞，本市新聞，本市新聞...

本市新聞，本市新聞，本市新聞...

本市新聞，本市新聞，本市新聞...

本市新聞，本市新聞，本市新聞...

本市新聞，本市新聞，本市新聞...

本市新聞，本市新聞，本市新聞...

本市新聞，本市新聞，本市新聞...

本市新聞，本市新聞，本市新聞...

本市新聞，本市新聞，本市新聞...

應有懸國旗之自由

應有懸國旗之自由，應有懸國旗之自由，應有懸國旗之自由...

應有懸國旗之自由，應有懸國旗之自由，應有懸國旗之自由...

應有懸國旗之自由，應有懸國旗之自由，應有懸國旗之自由...

應有懸國旗之自由，應有懸國旗之自由，應有懸國旗之自由...

應有懸國旗之自由，應有懸國旗之自由，應有懸國旗之自由...

應有懸國旗之自由，應有懸國旗之自由，應有懸國旗之自由...

應有懸國旗之自由，應有懸國旗之自由，應有懸國旗之自由...

應有懸國旗之自由，應有懸國旗之自由，應有懸國旗之自由...

應有懸國旗之自由，應有懸國旗之自由，應有懸國旗之自由...

應有懸國旗之自由，應有懸國旗之自由，應有懸國旗之自由...

應有懸國旗之自由，應有懸國旗之自由，應有懸國旗之自由...

應有懸國旗之自由，應有懸國旗之自由，應有懸國旗之自由...

應有懸國旗之自由，應有懸國旗之自由，應有懸國旗之自由...

應有懸國旗之自由，應有懸國旗之自由，應有懸國旗之自由...

應有懸國旗之自由，應有懸國旗之自由，應有懸國旗之自由...

應有懸國旗之自由，應有懸國旗之自由，應有懸國旗之自由...

應有懸國旗之自由，應有懸國旗之自由，應有懸國旗之自由...

應有懸國旗之自由，應有懸國旗之自由，應有懸國旗之自由...

應有懸國旗之自由，應有懸國旗之自由，應有懸國旗之自由...

應有懸國旗之自由，應有懸國旗之自由，應有懸國旗之自由...

應有懸國旗之自由，應有懸國旗之自由，應有懸國旗之自由...

渝陷區公產

渝陷區公產，渝陷區公產，渝陷區公產...

渝陷區公產，渝陷區公產，渝陷區公產...

渝陷區公產，渝陷區公產，渝陷區公產...

渝陷區公產，渝陷區公產，渝陷區公產...

渝陷區公產，渝陷區公產，渝陷區公產...

渝陷區公產，渝陷區公產，渝陷區公產...

渝陷區公產，渝陷區公產，渝陷區公產...

渝陷區公產，渝陷區公產，渝陷區公產...

渝陷區公產，渝陷區公產，渝陷區公產...

渝陷區公產，渝陷區公產，渝陷區公產...

渝陷區公產，渝陷區公產，渝陷區公產...

渝陷區公產，渝陷區公產，渝陷區公產...

渝陷區公產，渝陷區公產，渝陷區公產...

渝陷區公產，渝陷區公產，渝陷區公產...

渝陷區公產，渝陷區公產，渝陷區公產...

渝陷區公產，渝陷區公產，渝陷區公產...

渝陷區公產，渝陷區公產，渝陷區公產...

渝陷區公產，渝陷區公產，渝陷區公產...

渝陷區公產，渝陷區公產，渝陷區公產...

渝陷區公產，渝陷區公產，渝陷區公產...

渝陷區公產，渝陷區公產，渝陷區公產...

本報復刊小啓

本報復刊小啓，本報復刊小啓，本報復刊小啓...

本報復刊小啓，本報復刊小啓，本報復刊小啓...

本報復刊小啓，本報復刊小啓，本報復刊小啓...

本報復刊小啓，本報復刊小啓，本報復刊小啓...

本報復刊小啓，本報復刊小啓，本報復刊小啓...

本報復刊小啓，本報復刊小啓，本報復刊小啓...

本報復刊小啓，本報復刊小啓，本報復刊小啓...

本報復刊小啓，本報復刊小啓，本報復刊小啓...

本報復刊小啓，本報復刊小啓，本報復刊小啓...

本報復刊小啓，本報復刊小啓，本報復刊小啓...

本報復刊小啓，本報復刊小啓，本報復刊小啓...

本報復刊小啓，本報復刊小啓，本報復刊小啓...

本報復刊小啓，本報復刊小啓，本報復刊小啓...

本報復刊小啓，本報復刊小啓，本報復刊小啓...

本報復刊小啓，本報復刊小啓，本報復刊小啓...

本報復刊小啓，本報復刊小啓，本報復刊小啓...

本報復刊小啓，本報復刊小啓，本報復刊小啓...

本報復刊小啓，本報復刊小啓，本報復刊小啓...

本報復刊小啓，本報復刊小啓，本報復刊小啓...

本報復刊小啓，本報復刊小啓，本報復刊小啓...

本報復刊小啓，本報復刊小啓，本報復刊小啓...

杏花酒樓 高貴禮堂 租費從廉

好美華學園 高貴禮堂 租費從廉

陽痿 治限日 早洩 治限日

痰敵 除痰止咳 化痰止咳

順風牌 汽水 汽水

品需必之期暑 汽水 汽水

報晶 The Crystal 發行 A.L. Teodoro



家專影攝流一第上海 館相照術藝泰國

日三月八年國民於新報本 庫四分貳幣國幣貳日今 五期星 日一月四年捌拾貳國民華中

各路華軍大舉反攻 克復石龍包圍南昌 贛南高安發生激烈爭奪戰

地圖 顯示了中國南方的軍事行動區域，包括石龍、南昌、高安等地。

抱定抗戰決心 中國現決不談和 英當局亦聲明未提議調解

本報具有廿一年的歷史是小型報紙其批評時局普及全國刊登廣告最有力

法界懸旗糾紛

今日可復業

被捕市民昨日均已釋放

【本報訊】法界各界昨日（廿日）下午，在法界各界代表團與法租界工部局代表團之談判中，已告一段落。據悉：法界各界代表團，對於法租界工部局代表團所提之條件，已表示原則上之同意。惟對於其中若干細節，尚須進一步之商榷。據法界代表團發言人稱：談判已告一段落，被捕之市民，昨日均已釋放。法界各界代表團，將於今日（廿一日）下午，在法界各界代表團與法租界工部局代表團之談判中，正式簽署復業協議。據悉：復業協議之簽署，將使法界之秩序，得以恢復。法界各界代表團，將於今日（廿一日）下午，在法界各界代表團與法租界工部局代表團之談判中，正式簽署復業協議。據悉：復業協議之簽署，將使法界之秩序，得以恢復。

交涉經過

法界各界代表團，於昨日（廿日）下午，在法界各界代表團與法租界工部局代表團之談判中，已告一段落。據悉：法界各界代表團，對於法租界工部局代表團所提之條件，已表示原則上之同意。惟對於其中若干細節，尚須進一步之商榷。據法界代表團發言人稱：談判已告一段落，被捕之市民，昨日均已釋放。法界各界代表團，將於今日（廿一日）下午，在法界各界代表團與法租界工部局代表團之談判中，正式簽署復業協議。據悉：復業協議之簽署，將使法界之秩序，得以恢復。

河內

【本報訊】河內各界，昨日（廿日）下午，在河內各界代表團與法租界工部局代表團之談判中，已告一段落。據悉：河內各界代表團，對於法租界工部局代表團所提之條件，已表示原則上之同意。惟對於其中若干細節，尚須進一步之商榷。據河內代表團發言人稱：談判已告一段落，被捕之市民，昨日均已釋放。河內各界代表團，將於今日（廿一日）下午，在河內各界代表團與法租界工部局代表團之談判中，正式簽署復業協議。據悉：復業協議之簽署，將使河內之秩序，得以恢復。

租界已通

【本報訊】法界租界，昨日（廿日）下午，在法界租界代表團與法租界工部局代表團之談判中，已告一段落。據悉：法界租界代表團，對於法租界工部局代表團所提之條件，已表示原則上之同意。惟對於其中若干細節，尚須進一步之商榷。據法界租界代表團發言人稱：談判已告一段落，被捕之市民，昨日均已釋放。法界租界代表團，將於今日（廿一日）下午，在法界租界代表團與法租界工部局代表團之談判中，正式簽署復業協議。據悉：復業協議之簽署，將使法界租界之秩序，得以恢復。

外安界

【本報訊】法界外安界，昨日（廿日）下午，在法界外安界代表團與法租界工部局代表團之談判中，已告一段落。據悉：法界外安界代表團，對於法租界工部局代表團所提之條件，已表示原則上之同意。惟對於其中若干細節，尚須進一步之商榷。據法界外安界代表團發言人稱：談判已告一段落，被捕之市民，昨日均已釋放。法界外安界代表團，將於今日（廿一日）下午，在法界外安界代表團與法租界工部局代表團之談判中，正式簽署復業協議。據悉：復業協議之簽署，將使法界外安界之秩序，得以恢復。

本市

【本報訊】本市各界，昨日（廿日）下午，在法租界工部局代表團與法租界工部局代表團之談判中，已告一段落。據悉：本市各界代表團，對於法租界工部局代表團所提之條件，已表示原則上之同意。惟對於其中若干細節，尚須進一步之商榷。據本市代表團發言人稱：談判已告一段落，被捕之市民，昨日均已釋放。本市各界代表團，將於今日（廿一日）下午，在法租界工部局代表團與法租界工部局代表團之談判中，正式簽署復業協議。據悉：復業協議之簽署，將使本市之秩序，得以恢復。

法界

【本報訊】法界各界，昨日（廿日）下午，在法界各界代表團與法租界工部局代表團之談判中，已告一段落。據悉：法界各界代表團，對於法租界工部局代表團所提之條件，已表示原則上之同意。惟對於其中若干細節，尚須進一步之商榷。據法界各界代表團發言人稱：談判已告一段落，被捕之市民，昨日均已釋放。法界各界代表團，將於今日（廿一日）下午，在法界各界代表團與法租界工部局代表團之談判中，正式簽署復業協議。據悉：復業協議之簽署，將使法界之秩序，得以恢復。

應有懸國旗之自由

【本報訊】法租界各界，昨日（廿日）下午，在法租界各界代表團與法租界工部局代表團之談判中，已告一段落。據悉：法租界各界代表團，對於法租界工部局代表團所提之條件，已表示原則上之同意。惟對於其中若干細節，尚須進一步之商榷。據法租界各界代表團發言人稱：談判已告一段落，被捕之市民，昨日均已釋放。法租界各界代表團，將於今日（廿一日）下午，在法租界各界代表團與法租界工部局代表團之談判中，正式簽署復業協議。據悉：復業協議之簽署，將使法租界之秩序，得以恢復。

不堪勒索停業

【本報訊】法租界各界，昨日（廿日）下午，在法租界各界代表團與法租界工部局代表團之談判中，已告一段落。據悉：法租界各界代表團，對於法租界工部局代表團所提之條件，已表示原則上之同意。惟對於其中若干細節，尚須進一步之商榷。據法租界各界代表團發言人稱：談判已告一段落，被捕之市民，昨日均已釋放。法租界各界代表團，將於今日（廿一日）下午，在法租界各界代表團與法租界工部局代表團之談判中，正式簽署復業協議。據悉：復業協議之簽署，將使法租界之秩序，得以恢復。

租界已通

【本報訊】法界租界，昨日（廿日）下午，在法界租界代表團與法租界工部局代表團之談判中，已告一段落。據悉：法界租界代表團，對於法租界工部局代表團所提之條件，已表示原則上之同意。惟對於其中若干細節，尚須進一步之商榷。據法界租界代表團發言人稱：談判已告一段落，被捕之市民，昨日均已釋放。法界租界代表團，將於今日（廿一日）下午，在法界租界代表團與法租界工部局代表團之談判中，正式簽署復業協議。據悉：復業協議之簽署，將使法界租界之秩序，得以恢復。

外安界

【本報訊】法界外安界，昨日（廿日）下午，在法界外安界代表團與法租界工部局代表團之談判中，已告一段落。據悉：法界外安界代表團，對於法租界工部局代表團所提之條件，已表示原則上之同意。惟對於其中若干細節，尚須進一步之商榷。據法界外安界代表團發言人稱：談判已告一段落，被捕之市民，昨日均已釋放。法界外安界代表團，將於今日（廿一日）下午，在法界外安界代表團與法租界工部局代表團之談判中，正式簽署復業協議。據悉：復業協議之簽署，將使法界外安界之秩序，得以恢復。

本市

【本報訊】本市各界，昨日（廿日）下午，在法租界工部局代表團與法租界工部局代表團之談判中，已告一段落。據悉：本市各界代表團，對於法租界工部局代表團所提之條件，已表示原則上之同意。惟對於其中若干細節，尚須進一步之商榷。據本市代表團發言人稱：談判已告一段落，被捕之市民，昨日均已釋放。本市各界代表團，將於今日（廿一日）下午，在法租界工部局代表團與法租界工部局代表團之談判中，正式簽署復業協議。據悉：復業協議之簽署，將使本市之秩序，得以恢復。

法界

【本報訊】法界各界，昨日（廿日）下午，在法界各界代表團與法租界工部局代表團之談判中，已告一段落。據悉：法界各界代表團，對於法租界工部局代表團所提之條件，已表示原則上之同意。惟對於其中若干細節，尚須進一步之商榷。據法界各界代表團發言人稱：談判已告一段落，被捕之市民，昨日均已釋放。法界各界代表團，將於今日（廿一日）下午，在法界各界代表團與法租界工部局代表團之談判中，正式簽署復業協議。據悉：復業協議之簽署，將使法界之秩序，得以恢復。

高街粵茶 大小適宜 隨意小酌 租費從廉

杏花酒樓

華貴雅堂

好美真學園

妙法

妙法

陽痿早洩 可以早洩 治癒日期 保期如久

早洩

報晶

The Crystal

發行人 A.L. Teodoro

八二號三九路 號八九一第口路法租

新報之益得自無窮即如欲訂閱者請向本報

《號五十九百七千三期》

家專影攝流一第上海

館相照術藝泰國

化裝相照化演新製定

機良影攝光春此際

號五五路京州法地

號三九三〇九路電

日三月三年八國民於創報本 廉印分氣海通集張日今 五期星 日一廿月四年兩拾貳國民華中

各路華軍大舉反攻

克復石龍包圍南昌

贛南高安發生激烈爭奪戰

綏省重要據點多處已收復



【本報訊】各路華軍，自四月九日起，在贛南一帶，大舉反攻。據悉：華軍在石龍一帶，與敵發生激烈爭奪戰。目前，石龍已告克復。此外，華軍在南昌一帶，亦正積極包圍中。據悉：華軍在贛南一帶，大舉反攻。目前，石龍已告克復。此外，華軍在南昌一帶，亦正積極包圍中。據悉：華軍在贛南一帶，大舉反攻。目前，石龍已告克復。此外，華軍在南昌一帶，亦正積極包圍中。

抱定抗戰決心

中國現決不談和

【本報訊】中國政府，抱定抗戰決心，現決不談和。據悉：中國政府，抱定抗戰決心，現決不談和。據悉：中國政府，抱定抗戰決心，現決不談和。據悉：中國政府，抱定抗戰決心，現決不談和。

北平華軍得手

【本報訊】北平華軍，近日在軍事行動中，取得重大勝利。據悉：北平華軍，近日在軍事行動中，取得重大勝利。據悉：北平華軍，近日在軍事行動中，取得重大勝利。據悉：北平華軍，近日在軍事行動中，取得重大勝利。

本報發行小冊

【本報訊】本報為擴大宣傳，特發行小冊。據悉：本報為擴大宣傳，特發行小冊。據悉：本報為擴大宣傳，特發行小冊。據悉：本報為擴大宣傳，特發行小冊。

本報發行小冊

品雷必之期暑

點特大即 汽鮮順

美利華 汽水 牌

品出司公水汽華美

品雷必之期暑

點特大即 汽鮮順

美利華 汽水 牌

品出司公水汽華美

Wege zum Volltext:

Binnenstrukturen

HAASDZ 2018

法租界懸旗糾紛

公文上海日
國駐滬總
特生夫婦
在卽、特於
(九)、在
舉行茶會
埠之比僑
界國民精神總動員協會、於本
二十二日、正式成立後、即電呈
及工作狀況、頃該會接奉全國
國民精神總動員會張秘書長

於昨晚十
有被捕者
有沒收者
三)此後
轉陳總
各代界
人、准
納稅會領
由應勿加
範圍之內
應直接向
法租界領
事提出請
求

滬市工務林立、
戰後人口突增、工
業更見
原料
或游
股市
稻草
上之

報載、比
領事古特
以離滬在
昨日(十九
、在
其署內、
招待本埠
本市各界
國民精神
總動員會
蔣會長、
報告成
立經過、

國民精神總動員會張秘書長
岳軍復電、慰勉有加、茲照
錄如次(上略)：國民精神
總動員、確屬抗戰建國之重
要方略、滬市淪陷已久、敵
寇誘脅甚甚、策動愛國之精
神、爲抗戰之聲援、尤屬必
要、貴會領導羣衆、組織協
會、並利用各界固有之職業
團體、普及分途進行之持久機關、深謀
碩劃、殊堪嘉許、應准備案、仍希繼續
努力、奉行不懈、是所厚望、張章叩

四行軍團長謝晉元、昨
出巡市內、各界熱烈歡迎、
昨日爲上海市農工商
各界響應國民政府精
神總動員運動、租界
出與、一致懸掛紀念、愛受
法租界當局取締、發生糾紛
、應由代表一國之尊嚴
、任何國民、均應懸掛具國
旗之自由與權利、同胞乎、
若等英勇果敢、愛護國家、
擁護最高領袖之行為、世人
同深欽仰、余前次發表「一
平來觀感及對時局感想」文
曾言：今後國民唯有本「國
本興亡、匹夫有責之義、各
守崗位、竭盡天職、以最高
立分支行、逐漸推
行至新舊、以期增
加後方生產、充實
抗戰力量云、

前今、又有形同
稅卡之機關出現於
後、致無款以應、
相率停業、致稻草
來源驟減、各工廠
受缺貨影響、對動
索陋規深表憤慨、
同時車船亦反對絡
索、一致反對云

釋放經過

統計、惟捕房消息、曾有被捕及華捕
各一人受傷、被捕者、經納稅會交涉
後、已於昨日下午四時許、陸續釋放、
聞在捕房拘押時、多數均忍痛絕食、愛
國熱情、充分流露、惟被捕者中有一
名、捕房以其呼口號稱爲僑胞、初不允
釋放、經再三交涉、亦允釋放、上海市
銀行業同業公會、對於會員銀行之遭受
騷擾干涉、一面致電對面、一面勸其先
行復業、各銀行感以公會誠懇之勸慰、
亦即照常營業、又黃金大戲院執照被吊
銷後、聞亦發還、定即日起照常營業云

報載、比
領事古特
以離滬在
昨日(十九
、在
其署內、
招待本埠
本市各界
國民精神
總動員會
蔣會長、
報告成
立經過、

國民精神總動員會張秘書長
岳軍復電、慰勉有加、茲照
錄如次(上略)：國民精神
總動員、確屬抗戰建國之重
要方略、滬市淪陷已久、敵
寇誘脅甚甚、策動愛國之精
神、爲抗戰之聲援、尤屬必
要、貴會領導羣衆、組織協
會、並利用各界固有之職業
團體、普及分途進行之持久機關、深謀
碩劃、殊堪嘉許、應准備案、仍希繼續
努力、奉行不懈、是所厚望、張章叩

四行軍團長謝晉元、昨
出巡市內、各界熱烈歡迎、
昨日爲上海市農工商
各界響應國民政府精
神總動員運動、租界
出與、一致懸掛紀念、愛受
法租界當局取締、發生糾紛
、應由代表一國之尊嚴
、任何國民、均應懸掛具國
旗之自由與權利、同胞乎、
若等英勇果敢、愛護國家、
擁護最高領袖之行為、世人
同深欽仰、余前次發表「一
平來觀感及對時局感想」文
曾言：今後國民唯有本「國
本興亡、匹夫有責之義、各
守崗位、竭盡天職、以最高
立分支行、逐漸推
行至新舊、以期增
加後方生產、充實
抗戰力量云、

前今、又有形同
稅卡之機關出現於
後、致無款以應、
相率停業、致稻草
來源驟減、各工廠
受缺貨影響、對動
索陋規深表憤慨、
同時車船亦反對絡
索、一致反對云

消息聞
與西捕同
日法文上
開收此
聞聞開
浦江
政府
、以

法租界懸旗糾紛

公文上海日
國駐滬總
特生夫婦
在卽、特於
(九)、在
舉行茶會
埠之比僑
界國民精神總動員協會、於本
二十二日、正式成立後、即電呈
及工作狀況、頃該會接奉全國
國民精神總動員會張秘書長

滬市工務林立、
戰後人口突增、工
業更見
原料
或游
股市
稻草
上之

報載、比
領事古特
以離滬在
昨日(十九
、在
其署內、
招待本埠
本市各界
國民精神
總動員會
蔣會長、
報告成
立經過、

國民精神總動員會張秘書長
岳軍復電、慰勉有加、茲照
錄如次(上略)：國民精神
總動員、確屬抗戰建國之重
要方略、滬市淪陷已久、敵
寇誘脅甚甚、策動愛國之精
神、爲抗戰之聲援、尤屬必
要、貴會領導羣衆、組織協
會、並利用各界固有之職業
團體、普及分途進行之持久機關、深謀
碩劃、殊堪嘉許、應准備案、仍希繼續
努力、奉行不懈、是所厚望、張章叩

四行軍團長謝晉元、昨
出巡市內、各界熱烈歡迎、
昨日爲上海市農工商
各界響應國民政府精
神總動員運動、租界
出與、一致懸掛紀念、愛受
法租界當局取締、發生糾紛
、應由代表一國之尊嚴
、任何國民、均應懸掛具國
旗之自由與權利、同胞乎、
若等英勇果敢、愛護國家、
擁護最高領袖之行為、世人
同深欽仰、余前次發表「一
平來觀感及對時局感想」文
曾言：今後國民唯有本「國
本興亡、匹夫有責之義、各
守崗位、竭盡天職、以最高
立分支行、逐漸推
行至新舊、以期增
加後方生產、充實
抗戰力量云、

前今、又有形同
稅卡之機關出現於
後、致無款以應、
相率停業、致稻草
來源驟減、各工廠
受缺貨影響、對動
索陋規深表憤慨、
同時車船亦反對絡
索、一致反對云

法租界懸旗糾紛

公文上海日
國駐滬總
特生夫婦
在卽、特於
(九)、在
舉行茶會
埠之比僑
界國民精神總動員協會、於本
二十二日、正式成立後、即電呈
及工作狀況、頃該會接奉全國
國民精神總動員會張秘書長

滬市工務林立、
戰後人口突增、工
業更見
原料
或游
股市
稻草
上之

報載、比
領事古特
以離滬在
昨日(十九
、在
其署內、
招待本埠
本市各界
國民精神
總動員會
蔣會長、
報告成
立經過、

國民精神總動員會張秘書長
岳軍復電、慰勉有加、茲照
錄如次(上略)：國民精神
總動員、確屬抗戰建國之重
要方略、滬市淪陷已久、敵
寇誘脅甚甚、策動愛國之精
神、爲抗戰之聲援、尤屬必
要、貴會領導羣衆、組織協
會、並利用各界固有之職業
團體、普及分途進行之持久機關、深謀
碩劃、殊堪嘉許、應准備案、仍希繼續
努力、奉行不懈、是所厚望、張章叩

四行軍團長謝晉元、昨
出巡市內、各界熱烈歡迎、
昨日爲上海市農工商
各界響應國民政府精
神總動員運動、租界
出與、一致懸掛紀念、愛受
法租界當局取締、發生糾紛
、應由代表一國之尊嚴
、任何國民、均應懸掛具國
旗之自由與權利、同胞乎、
若等英勇果敢、愛護國家、
擁護最高領袖之行為、世人
同深欽仰、余前次發表「一
平來觀感及對時局感想」文
曾言：今後國民唯有本「國
本興亡、匹夫有責之義、各
守崗位、竭盡天職、以最高
立分支行、逐漸推
行至新舊、以期增
加後方生產、充實
抗戰力量云、

前今、又有形同
稅卡之機關出現於
後、致無款以應、
相率停業、致稻草
來源驟減、各工廠
受缺貨影響、對動
索陋規深表憤慨、
同時車船亦反對絡
索、一致反對云

四里路

杏花酒樓

大小筵席 隨意小酌
電話 五五三九
九四六二九

高尚粵菜

華貴禮堂 租費從廉

美好樂園

高尚 娛樂 妙演劇團
電話 二〇八七 號四〇二一 路

各路華軍大舉反攻

紙鼻祖銷路普

應有懸國旗之自

謝晉元團長昨發告同胞
四行軍團長謝晉元、昨
出巡市內、各界熱烈歡迎、
昨日爲上海市農工商
各界響應國民政府精
神總動員運動、租界
出與、一致懸掛紀念、愛受
法租界當局取締、發生糾紛
、應由代表一國之尊嚴
、任何國民、均應懸掛具國
旗之自由與權利、同胞乎、
若等英勇果敢、愛護國家、
擁護最高領袖之行為、世人
同深欽仰、余前次發表「一
平來觀感及對時局感想」文
曾言：今後國民唯有本「國
本興亡、匹夫有責之義、各
守崗位、竭盡天職、以最高
立分支行、逐漸推
行至新舊、以期增
加後方生產、充實
抗戰力量云、

日捕加薪

與西捕同待遇二十
日法文上海日報載
聞聞開
浦江
政府
、以

滑稽一重

陰兵們到
旁邊一個年紀大一點的、陪
都已經弄得很亂了。他們的
一來、吃不慣天上的東西、
飯匙毒藥在內、所以無論什
帶到了天上、燃料就無處
是他們把人家桌子椅子
中什物、有跟得幾、他們
還說有許多想講、就叫我們

本報具有廿一年的歷史是小型報紙鼻祖銷路普

法租界懸旗糾紛

上海市商會及法租界各商店、應請各商店、即由法租界巡捕房、分別發出緊急公函、(一)納稅會決辦法、此為顧念中法邦交、應請大、而利進行、全部復業云、

交涉經過

於昨晚十時、往法租界巡捕房、徐陳述商民意見外、並要、有被捕諸人、應請即予釋放、有沒收旗、應由捕房送周、(三)此後凡遇紀念節日、仍、備助總監當表示、節日、各代表再往、當得辦法(一)人、准即釋放、(二)在捕納稅會備函領回、(三)至由應放勿加干涉、(四)因不範圍之內、應直接向法總領

釋放經過

統計、惟捕房消息、曾有各一人受傷、被捕諸人、經後、已於昨日下午四時許、陸續、開在捕房拘押時、多數均拒捕、愛國熱情、充分流露、惟被捕諸人中、有二名、捕房以其呼口號稱爲煽動、初不允釋放、經再三交涉、亦允釋放、上海市銀行業同業公會、對於會員銀行之遭受懸旗干涉、一面致意對會、一面勸其先、亦即照常營業、又黃金大戲院執照被吊銷後、聞亦發還、定即日起照常營業云

法租界糾紛

狂業 爲奸 物、而工 之農產品 來自內地 城、致一 斷居奇、 亦工業 少者、故 近忽有 業公會董 (一名翟 南碼頭開

公共租界 日捕加薪 與西捕同待遇二十 日法文上海日報載 公共租界巡捕 房之日籍巡捕、 將增加薪水百分之 三十五、而與西捕 滑稽 小說 一重

有懸國旗之自 謝元團長昨發告同胞 領袖之主張爲主 意志、而腳踏實 國、孤島之同胞 年開始、隨人人 深望孤島同胞、 效祖國之責任、 有力出力、以完 之基礎、余今 同孤島同胞呼 受之苦痛、區 之黑暗、一切惡 華民族偉大光榮 鴻、切望孤島同 有責、出錢出力 負責、負責、望 前、又有形同 稅卡之機關出現於 後、致無款以應、 相率停業、致租車 當升 公共租界 日捕加薪 與西捕同待遇二十 日法文上海日報載 公共租界巡捕 房之日籍巡捕、 將增加薪水百分之 三十五、而與西捕 滑稽 小說 一重

四里路 杏花園 酒樓 電話 九四六二九

華貴禮堂 租費從廉 大小筵席 隨意小酌

美好 娛樂 妙演真出色 電話 二〇〇八七號四〇二一

各路華軍大舉反攻

本報具有廿一年的歷史是小型報紙鼻祖銷路普

各商店遵勸

今日可復業

被捕市民昨日均已釋放

關於以後自由懸旗一節解待解決
上海市商會及法租界納稅華人會、爲勸導法租界商店復業、於昨晨、分別發出緊急通告云、(一)市商會、(二)法租界交際會、此次懸旗糾紛、係由總捕房執行職務、過分所致、各界引爲遺憾、除電請政府交涉外、應請各商店、即先復業、以堅毅貞固之態度、靜候合法合理之解決、(三)納稅會、(四)查關於懸旗事件、業經本會推員向法租界商會商解、決辦法、此爲應辦大、而利於全部復業云、

交涉經過

於昨晚十時、有被捕諸人、有沒收懸旗、(三)此後凡遇轉陳總監後、各代表再往、人、准即釋放、納稅會備函領、由懸旗勿加干、範圍之內、應

釋放經過

統計、惟捕房各一人受傷、後、已於昨日、開在捕房拘押、國熱情、充分、名、捕房以其、釋放、經再三、銀行業同業公、懸旗干涉、行復業、各報、亦即照常營業、銷後、聞亦發

法報記

發生、外間之、日因大規模懸、質之表示、故、當時並無

報紙鼻祖銷

任何國國民

應有懸國旗

謝晉元團長昨

四行孤軍團長謝晉元、昨發出爲市籍懸掛國旗敬告同胞書、大要如下、昨日爲上海市農工商各界響應國民政府精

西北金融網

已完竣

中國銀行、設蘭州分行、特設蘭州爲經理、已抵蘭籌備一切、據談際此非常時期、中央爲完成西北金融之計劃、特命

任何國國民

應有懸國旗之自

謝晉元團長昨告同胞、四行孤軍團長謝晉元、昨發出爲市籍懸掛國旗敬告同胞書、大要如下、昨日爲上海市農工商各界響應國民政府精

捕加薪

捕房待遇二十、浦上每日報載、公共租界巡捕、加薪百分之五、而與西捕

一重

陰兵們到、一團年紀大一點的、弄得亂了。他們的、吃不慣天上的東西、藥在內、所以無論什、天的恩料、就不敢吃、就把人家的桌子椅子、有許多人想、就叫我們

四里路

杏花酒樓

電話 五五三九 九四六二九

專辦

高尚粵菜

大小筵席 隨意小酌

華貴禮堂 租費從廉

美好

高尚 好萊塢劇團 娛樂 妙演出色 之新 之溫柔鄉 團地 聯歌求樂

電話 二〇八七 號四〇二一 路

各路華軍大舉反攻

報晶

The Crystal

發行人 A. L. Teodoro

八二四三九路道 號九二路日漢法新

紙新之送寄妥包郵照按立號掛准特政郵華中

《號四廿字〇號記登部工和租共公》

《號五十九百七千三第》



家專影攝流一第上海
館相照術藝泰國

化族貴相照化濟經價定
機良影攝光春此際

號五路京南址地
號三九三〇九話電

日三月三年八國民於前報本 厘四分幣國幣張貳日今 五期星 日一廿月四年捌拾貳國民華中

本報志實報道新聞說人人要說的話文字淺近趣味濃厚是最平民化的刊物



各路華軍大舉反攻 克復石龍包圍南昌

贛南高安發生激烈爭奪戰 緩督重要據點多處已收復

【本報廿日專電】贛南各路華軍，自三月下旬起，即紛紛發動反攻。目前，高安、石龍、包圍南昌等處，均發生激烈爭奪戰。據悉，各路華軍已收復重要據點多處，並正向南昌推進中。

抱定抗戰決心 中國現決不談和

英當局亦聲明未提議調停

【本報廿日專電】中國政府抱定抗戰決心，現決不談和。英當局亦聲明，目前尚未提議調停。此舉顯示中英兩國在抗戰問題上立場一致，絕不向侵略者屈服。

蘇聯戰事 華軍甚得手

宿潼流潼各得勝 華軍主動

【本報廿日專電】蘇聯戰事進展順利，華軍在宿潼、流潼等處均取得勝利。華軍主動進攻，戰果輝煌，顯示其戰鬥力已顯著提高。

梅吳揚揚 梅吳揚揚

【本報廿日專電】梅吳揚揚，戰事進展順利。梅吳揚揚，戰事進展順利。梅吳揚揚，戰事進展順利。

本報復刊小啓

本報自一月廿一日起，由美商發行，恢復出版。本報宗旨不變，仍致力於報導抗戰實況，服務廣大讀者。特此聲明。

杏花酒樓

華貴禮堂 租費從廉

高貴粵菜 大小筵席 隨意小酌

路馬四里師 九六二九九

高崗野地樂園

娛樂 妙畫園 妙畫園

地址：高崗野地

陽痿早洩

可以早洩 治癒早洩

本藥專治陽痿早洩，功效顯著。各大藥房均有代售。

各商店遵勸 今日可復業

被捕市民昨日均已釋放

開放以後自由懸旗一節暫待解決

【本報訊】上海市政府及各機關，為維持社會秩序，並使各商店早日復業，特頒布各商店遵勸。凡各商店應於今日開始復業。對於被捕市民，昨日均已釋放。關於開放以後自由懸旗一節，暫待解決。

河內 航郵已通

【本報訊】河內航郵已通，交通恢復。市民可通過航空郵件寄送包裹，極大方便了市民生活。

不推銷索停業

【本報訊】不推銷索停業，市民應注意。政府呼籲市民停止推銷索，以維護社會秩序。

任何國民 應有懸國旗之自由

【本報訊】任何國民應有懸國旗之自由。政府重申此項權利，以彰顯國家尊嚴。

淪陷區公產

【本報訊】淪陷區公產，政府正積極處理。將採取措施保護公產，防止流失。

小滑稽一重天

【本報訊】小滑稽一重天，笑話笑話。本報特輯笑話，供市民一笑。

陽痿早洩 治癒早洩

本藥專治陽痿早洩，功效顯著。各大藥房均有代售。

夏季必備之品

順風牌 汽水 汽水

夏季必備之品，順風牌汽水，清爽解渴。各大超市均有代售。

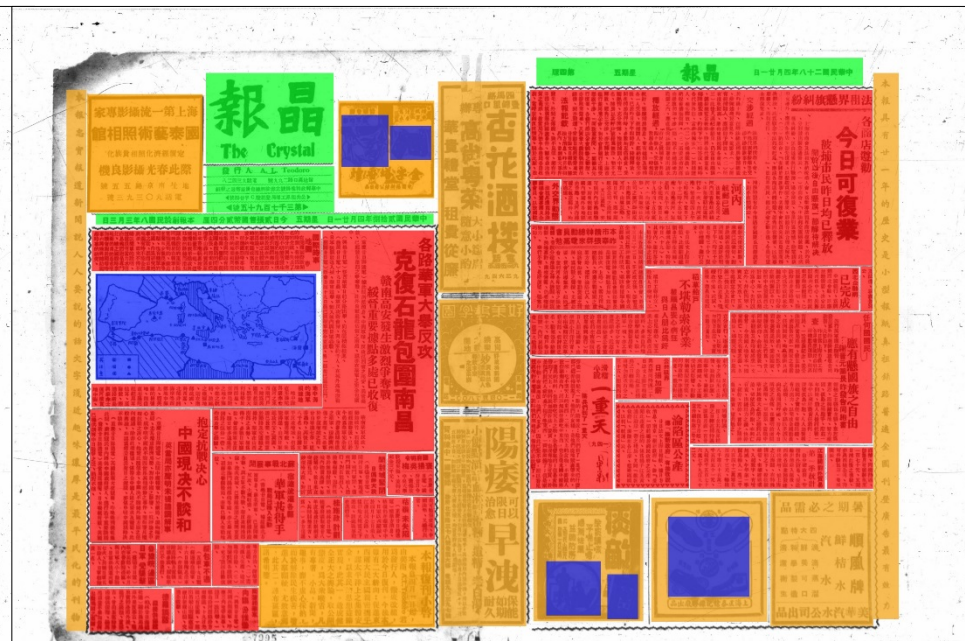
Wege zum Volltext:

Manuelle Segmentierung (Crowd)

JADH 2018

Segmentation - I

- Page segmentation (pattern recognition/computer vision)
 - Analyze layout of page, use page-internal structures
 - Identify semantic units
 - Generate co-ordinates, relate them to items, store in DB



Segmentation - II

- Page segmentation (crowdsourcing)
 - Pilot project with Pallas Ludens GmbH
 - Let the crowd help analyzing the pages
 - Identify and label four item types:
 - image/drawing
 - article
 - advertisement
 - additional information
 - Supervised
 - Non-Chinese speaking community!



晶報


本報每日出版除星期日及國慶日外全年無間
 零售每份五分
 本埠每月一元二角
 外埠每月一元五角
 廣告刊例
 第一版每行一元
 第二版每行八角
 第三版每行六角
 第四版每行四角
 長期廣告另議

定報
 本報每月一元二角
 外埠每月一元五角
 廣告刊例
 第一版每行一元
 第二版每行八角
 第三版每行六角
 第四版每行四角
 長期廣告另議

科發白濁丸

治淋病特效藥
 此丸專治男女白濁、赤白帶下、腰酸背痛、小便頻數、淋漓不盡、遺精早泄、婦女經閉、赤白帶下、子宮炎、附件炎、一切淋病、服此丸無不立效。每盒一元。

請吸新出 長城牌國貨香煙



長城牌香煙，品質優良，口味醇厚，為國貨之光。每包十支，售價一元。

上海五洲大藥房

自來血、海狗丸、樹皮丸
 自來血：補血強身，治貧血、頭暈、眼花、心悸、失眠、月經不調等症。每瓶二元。
 海狗丸：補腎壯陽，治腰酸背痛、陽痿早洩、遺精滑精、精神不振等症。每瓶二元。
 樹皮丸：止咳化痰，治傷風感冒、咳嗽氣喘、痰多胸悶等症。每瓶二元。

德六零六

梅毒特效藥
 德國六零六，專治梅毒、淋病、下疳、橫痃、魚口、便毒、濕疹、皮膚瘙癢等症。每瓶一元。

張世楷西醫

專治各種疑難雜症
 張世楷醫師，醫學精湛，經驗豐富。專治各種疑難雜症，如：神經衰弱、失眠、頭痛、胃病、肝病、腎病、婦科疾病等。診所地址：上海南京路。

中法儲蓄會

儲蓄致富，積少成多
 中法儲蓄會，提供多種儲蓄方案，如：零存整取、整存整取、活期儲蓄等。手續簡便，利息優厚。歡迎各界人士參加。

萬應消癰丸

治一切瘡癤腫毒
 萬應消癰丸，專治一切瘡癤、腫毒、疔瘡、癰疽、乳癰、痔瘡、脫肛、瘻管、一切無名腫毒。每盒一元。

秋令咳嗽

止咳化痰，清肺潤燥
 秋令咳嗽，多因肺燥、氣逆所致。本藥專治咳嗽、氣喘、痰多、胸悶、肺癆、吐血、咯血、一切肺病。每瓶一元。

預防花柳奇藥 愛克憐

預防花柳，保護健康
 愛克憐，預防花柳、梅毒、淋病、下疳、橫痃、魚口、便毒、濕疹、皮膚瘙癢等症。每瓶一元。

萬隆老棧

馳名中外，貨真價實
 萬隆老棧，經營各種名產、土貨、雜貨。貨真價實，童叟無欺。歡迎各界人士光臨。

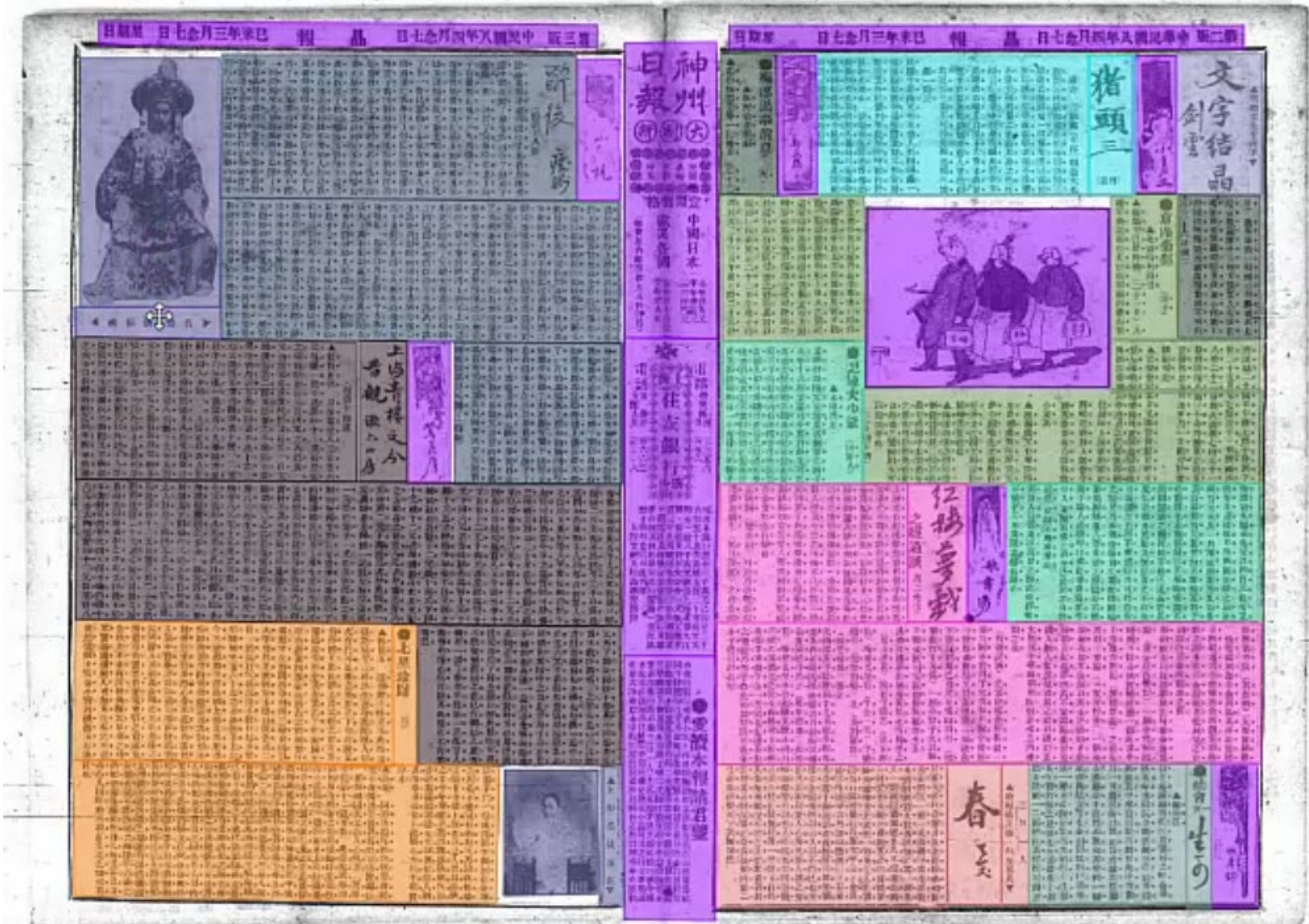
安氏補藥片

補血氣，助消化，強體力
 安氏補藥片，補血氣、助消化、強體力。專治貧血、頭暈、眼花、心悸、失眠、月經不調、食慾不振、消化不良、便秘、腹瀉、一切虛弱症。每瓶一元。



寒雲主人好古知書深探三代漢魏之神髓...
 本報每日出版除星期日及國慶日外全年無間
 零售每份五分
 本埠每月一元二角
 外埠每月一元五角
 廣告刊例
 第一版每行一元
 第二版每行八角
 第三版每行六角
 第四版每行四角
 長期廣告另議

Grouping semantic units



Outcome of segmentation pilot

1. Page segmentation can be outsourced to expert crowd
 - Requires supervision
 - Advanced user interfaces (high usability, efficiency)
 - Crowd should read Chinese (semantic grouping)
2. *Jingbao* 晶報 1919-21 completely segmented with qualified boxes, issues of April 1919 with semantic units
3. Further processing:
 - Partnership with Computational Knowledge Lab (知識計算實驗室), Department of Engineering Science and Ocean Engineering, Taiwan National University, <http://www.cklab.org/>
 - Seeking additional partners for collaboration!

Wege zum Volltext:

Ground truth

Funding for Ground Truth

Zusätzliche Mittel über Field of Focus 3 (ExStrat Heidelberg), Juli-
Oktober 2019

Gemeinsam mit Duncan Paterson und 4 WiHi (typing)

- GroundTruth bounding boxes mit Labels
 - Koordinaten in json und web-annotation Format
 - Jingbao April 1939 und April 1938
 - Stand: 70 folds, 6335 shapes, Ø 90,5 shapes/fold
- GroundTruth Texte
 - Blind double-keying
 - Jingbao April 1939 - 10 issues, 40 folds
 - Stand: ca. 245.000 Zeichen, lokales XML Format, Ø 6100 Zeichen/fold (Bearbeitung ca. 10h/fold)

Annotation tool



mode:select

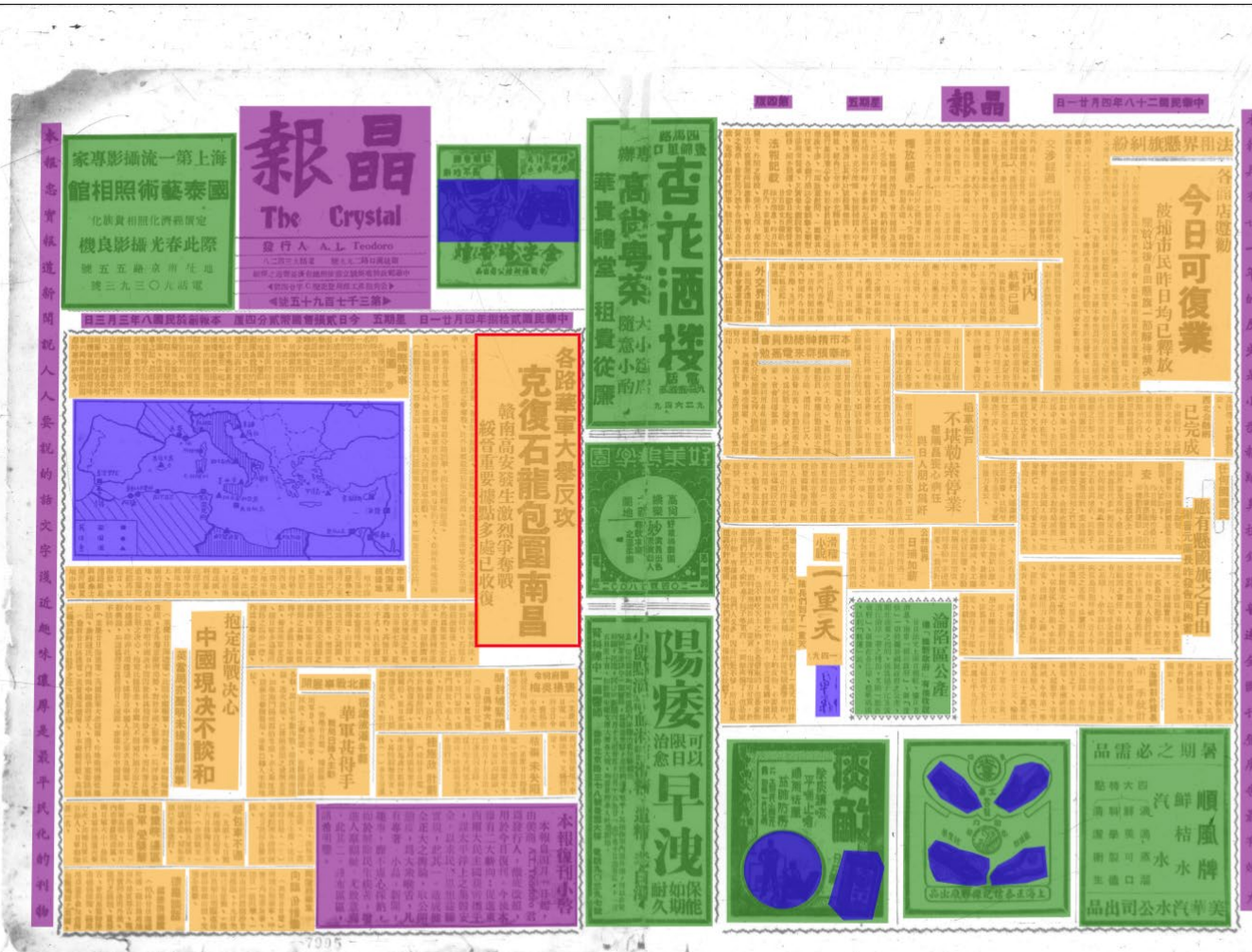


Image or Drawing

Article

Advertisement

Additional Information

SUBMIT



Entwicklung: eXist-Solutions

Annotationen in json

jb_3795_1939-04-21_0002+0003.json x

```
1 {
2   "items": [
3     {
4       "body": [
5         {
6           "value": {
7             "color": "purple",
8             "name": "additional",
9             "label": "Additional Information"
10          },
11          "type": "CategoryLabel"
12        }
13      ],
14      "created": "2019-07-10T16:14:26.772Z",
15      "target": [
16        {
17          "type": "SpecificResource",
18          "source": "https://kjc-sv002.kjc.uni-heidelberg.de:8080/fcgi-bin/iiprv.fcgi?IIIF=imageStorage/ecpo_new/jingbao/1939/04/jb_3795_1939-04-21_0002%252B0003.tif/full/full/0/default.jpg",
19          "id": "s-156277526676110",
20          "selector": {
21            "value": "<g transform=\"matrix(1 0 0 1 3462.22464 371.24256)\"><polygon points=\"-812.0176,-74.99264 812.0176,-74.99264 812.0176,74.99264 -812.0176,74.99264 \"/></g>",
22            "type": "SvgSelector"
23          }
24        }
25      ],
26      "type": "Annotation",
27      "id": "s-1562775266771"
28    },
29    {
30      "body": [
31        {
32          "value": {
33            "color": "purple",
34            "name": "additional",
35            "label": "Additional Information"
36          },
37          "type": "CategoryLabel"
38        }
39      ],
40      "created": "2019-07-10T16:14:48.631Z",
41      "target": [
42        {
43          "type": "SpecificResource",
44          "source": "https://kjc-sv002.kjc.uni-heidelberg.de:8080/fcgi-bin/iiprv.fcgi?IIIF=imageStorage/ecpo_new/jingbao/1939/04/jb_3795_1939-04-21_0002%252B0003.tif/full/full/0/default.jpg",
45          "id": "s-15627752886237",
46          "selector": {
47            "value": "<g transform=\"matrix(1 0 0 1 4331.59808 1719.37056)\"><polygon points=\"-26.18944,-1352.42848 53.00704,-1352.42848 26.23776,1352.42848 -53.00704,1352.42848 \"/></g>",
48            "type": "SvgSelector"
49          }
50        }
51      ],
52      "type": "Annotation",
53      "id": "s-1562775288631"
54    }
55  ]
56 }
```

Texte auf GitHub - 1

Lokales XML Schema

- `<fold>`
mit `@xml:id` und recto/verso-Paar
- `<div>` (division)
mit `page`, `article`, `image`, `advert`, `other`, `head` (for running head on top of page), `margin` (for margin)
- `<p>` (paragraph)
- `<lb/>` (line break) / `<pb/>` (page break)

Encoding “UTF-8”:

- Zeichen nach Vorlage, wo immer Unicode codepoint verfügbar
- Unlesebar: `&gaiji;` (“``“)
- Alle Schriftzeichen als double-space chars

Texte auf GitHub - 2

Besonderheit: Laufrichtung des Textes

`<div>` mit `@mode` und `@dir`

Default: `mode="vertical-rl"`

```
<div type='other' mode="vertical-rl">
  <div type='advert'>
    <!-- first contents -->
  </div>
  <div type='advert' mode="horizontal-tb" dir='ltr'>
    <!-- second contents -->
  </div>
  <div type='advert'>
    <!-- third contents -->
  </div>
</div>
```

<https://github.com/exc-asia-and-europe/ecpo>


```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <?xml-model href="https://raw.githubusercontent.com/exc-asia-and-europe/ecpo/master/ecpo.rng" type="application/x
3 <!DOCTYPE ecpo [
4     <!ENTITY gaiji "&#xe111;">
5 ]>
6 <fold xml:id="jb_3795_1939-04-21_0001+0004">
7     <!-- start recto page -->
8     <div type="page" page="recto" n="1" mode="vertical-rl" dir="rtl">
9         <div type="head" dir="rtl" mode="horizontal-tb">
10             <div>
11                 晶報<lb/>
12             </div>
13             <div dir="ltr" mode="horizontal-tb">
14                 The Crystal<lb/>
15                 發行人 A. L. Teodoro<lb/>
16             </div>
17             <div>
18                 館址漢口路二九九號 電話九三四二八<lb/>
19                 中華郵政特准掛號立券按照總包優益寄送之報紙<lb/>
20                 ◀公共租界工部局登記證C字廿四號▶<lb/>
21                 ▶第三千七百九十五號◀<lb/>
22                 中華民國貳拾捌年四月廿一日 星期五 今日兩全張售國幣貳分四厘 本報創刊於民國八年三月三日<lb/>
23             </div>
24         </div>
25         <div type="advert">
26             <div>
27                 煙中<lb/>
28                 鐵軍<lb/>
29                 清香<lb/>
30                 雋永<lb/>
31             </div>
32             <div mode="horizontal-tb" dir="rtl">
33                 &gaiji;明者&gaiji;<lb/>
34                 &gaiji;不吃虧<lb/>
35             </div>
36             <div type="image"></div>
37             <div type="image"></div>
```

Wege zum Volltext:

Machine learning

Impresso und dhSegment

DH2019 (Utrecht): Kontakt impresso

Folgegespräche insb. mit Sofia Oliveira

-> Einsatz dhSegment: “generic deep-learning framework for Historical Document Processing” <https://github.com/dhlab-epfl/dhSegment>

Ziele

- Unterstützung bei Document segmentation
- Trainieren der Modelle auf Unterscheidung 3 Hauptbereiche
 - Paratext (Kopfzeilen/running head und Marginalia)
 - Werbung und Bilder
 - Textbereiche
- Unterstützung der Maschine durch Entfernung der Bereiche mit Spezialfonts und Bildern (Abarbeitung einzeln)

Paratext

Erster Schritt

- Weitgehend regelmäßig
- Nur geringe Unterschiede im Text (zb. Seitenzahlen, Datum)
- Daten können für Metadaten genutzt werden
 - Editors, Publisher
 - Identifikation der Seiten
- Sonderfonts und Layout behindern Seiten- und Fonterkennung

Ziele

- Netzwerk trainieren auf Erkennung Header in ganz Jingbao (in Arbeit)
- Nutzung zum OCR Training (begrenzter Satz an Zeichen)
- Wenige unregelmäßige Zeichen und Fonts

Illustrationen und Werbung

- Werbung meist im Falz, teilweise schlecht erhalten
- Viele Wiederholungen von Ads innerhalb bestimmter Zeiträume
- Derzeit noch kein tracking von Ads zwischen verschiedenen Issues / Publikationen

Herausforderung:

- Zeichen/Fonts extrem variabel – Sonderzeichen, Spezialfonts, invertierte Zeichen, Text-in-Bild-in-Text, etc

Ziele:

- Modell zur Erkennung einzelner Ads
- Über Einzel-Fold hinaus
- Metadaten

Inhalte / Text

Ideal

- Gruppierung der Sammlung nach Fonts (historisch wurden nur eine sehr begrenzte Anzahl von CJK Fonts genutzt)
- Texte in „natural language“, legt Einsatz von Wörterbüchern für die Verarbeitung nahe

Aber

- historische Varianten (Zeichen, Wort) nicht in modernen Modellen
- Mischung von Schriften (Hant, Latn) und typische Wechsel in Laufrichtung führen zu Problemen in OCR
- GroundTruth Erstellung sehr zeitaufwändig, Verarbeitung mit extensivem pre-processing (comments, ornamentation, etc)

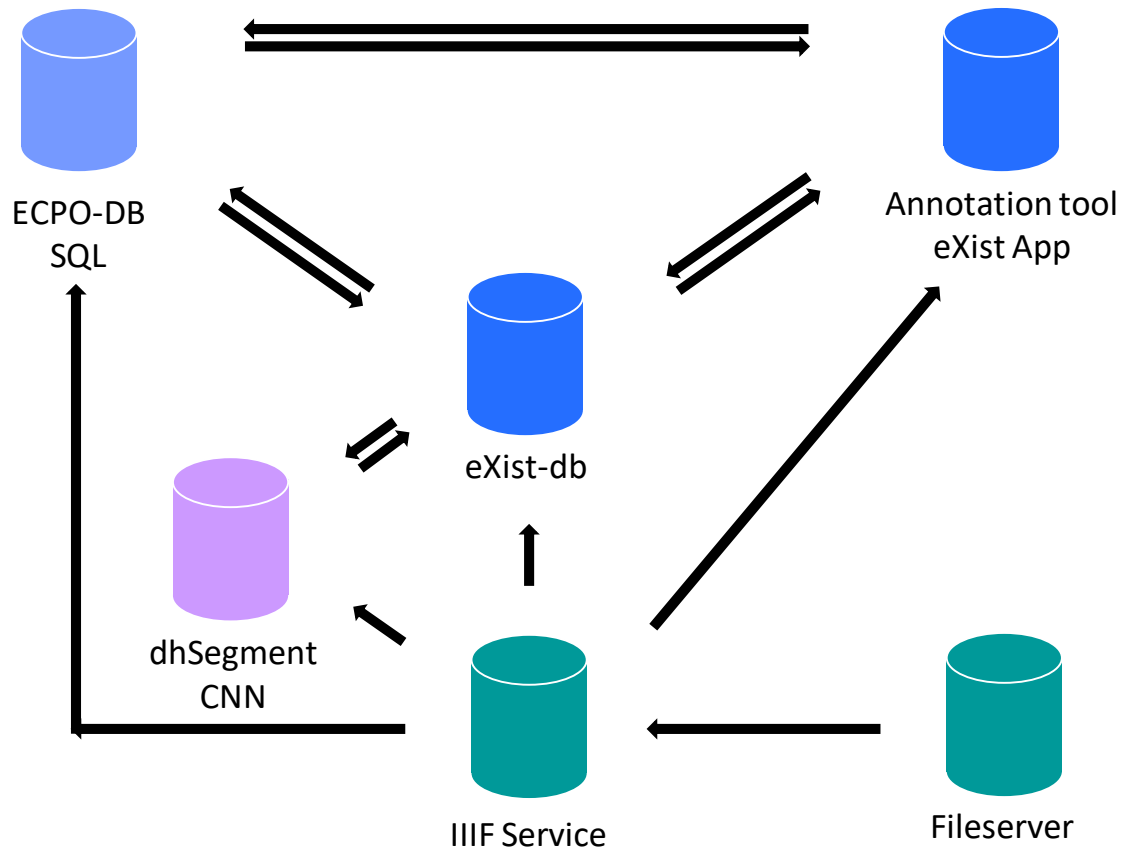
Ziele

- Verarbeitung einzelner kompletter Jahrgänge und Einbau in ECPO Search
- Ausreichende OCR Qualität für einfache Suchen, incl. Zeichenvarianten

Wege zum Volltext:

Nächste Schritte

ECPO – Erweiterte Infrastruktur



Herausforderungen

- Rechenkapazitäten unzureichend für Entwicklung (GW und GPU Server)
- Gewachsene Teilmodule müssen in neuer Infrastruktur integriert werden
 - Unit-ID's, file encodings, etc.
- Outsourced Module teilweise mangelhaft dokumentiert

Work in progress

- Generierung von .png aus Anno-Tool für CNN Training mit GroundTruth
- Verlinkung der Annotationen (manuell und Maschine) mit vorhandenen Annotationen in ECPO-SQL (Metadaten Artikel, Bild, Werbung)
- Stapelverarbeitung der Regionen im Scan nach Typ

Mittelfristig

- Modell auf Wechsel der Laufrichtung innerhalb semantischer Einheiten trainieren
- Modell für vertikale Zeilenerkennung trainieren (RtL, TtB)
- Erstellung Wörterbuch für Republikzeitliches geschriebenes Chinesisch

Langfristig

- Ausdehnung des Modells auf andere non-Latn script Materialien
- Verbesserung der CJK OCR für unterschiedliche hist. Perioden

Contact

Matthias Arnold

Heidelberg Centre for Transcultural Studies | HCTS
Karl Jaspers Centre
Voßstr. 2 | Building 4400 | Room 005b
69115 Heidelberg, Germany

Phone: +49 - 6221 - 54 4094
eMail: matthias.arnold@uni-hd.de
Web: <http://tinyurl.com/matthias-arnold>

References

Women and the Periodical Press in China's Global Twentieth Century: A Space of Their Own? Ed. by Joan Judge, Barbara Mittler and Michel Hockx, Cambridge University Press, 2018.

Arnold, Matthias, und Lena Hessel. „Transforming Data Silos into Knowledge: Early Chinese Periodicals Online (ECPO)“. Heidelberg: Heidelberg University Press, 2019. <https://doi.org/DOI:10.11588/heidok.00027325>.

Sung, Doris, Liying Sun, und Matthias Arnold. „The Birth of a Database of Historical Periodicals: Chinese Women's Magazines in the Late Qing and Early Republican Period“. *Tulsa Studies in Women's Literature* 33, Nr. 2 (2014): 227–37. <https://doi.org/10.1353/tsw.2014.0004>.

Oliveira, Sofia Ares, Benoit Seguin, and Frederic Kaplan. “dhSegment: A Generic Deep-Learning Approach for Document Segmentation.” *2018 16th International Conference on Frontiers in Handwriting Recognition (ICFHR)*, August 2018, 7–12. <https://doi.org/10.1109/ICFHR-2018.2018.00011>.