

Semantic segmentation and document layout recognition: Approaches towards full text recognition of early Chinese newspapers

Matthias Arnold | Uni Heidelberg · HCTS · HRA

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>

報晶

The Crystal

發行人 A. L. Teodoro

八二四號大馬路 號九二第口路地府

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
化裝貴相照化濟經價定
機良影攝光春此際
號五五路京州址地
號三九三〇九路電



專辦 高尙粵茶 大小適宜 隨意小酌 華貴禮堂 租費從廉

杏花酒樓

路馬四 口里第 九二第

日三月三年八國民於新華本 厘四分貳幣國幣張貳日今 五期星 日一廿月四年捌貳國民華中

本報自廿一年之歷史是小型報紙其粗銷路普通全國刊登廣告最有效力



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

（中央社訊）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。據悉，華軍在贛南一帶，正與敵展開激烈之拉鋸戰，高安、石龍等處，均為爭奪之焦點。同時，在綏遠省境內，華軍亦正積極反攻，多處重要據點，均已收復。此項反攻行動，係華軍自發動以來，最重大之進展，顯示華軍士氣高昂，戰鬥力強盛。

抱定抗戰決心 中國現狀

英當局亦擬 援手

（中央社訊）中國政府抱定抗戰決心，中國現狀。據悉，中國政府正積極籌劃抗戰事宜，並已獲得英政府之援助。英政府表示，將以各種方式，援助中國之抗戰事業。此項援助，將包括物資、技術及財政等方面。中國政府對此表示衷心感謝，並表示，將繼續堅持抗戰，直至取得最後之勝利。

各商店運動 今日可復業

被捕市民昨日均已釋放

關於自衛隊自由團被解散一事，業經政府與各團體代表，經過多次協商，業已達成協議。據悉，政府已決定，自衛隊自由團之解散，係屬行政命令，並非法律行為。因此，該團成員之被捕，均屬非法。政府已決定，將所有被捕之自衛隊自由團成員，一律釋放。同時，政府亦決定，將對該團成員之財產，予以保護。此項決定，已獲得各團體代表之同意。據悉，各商店亦將於今日恢復營業。此項決定，對於穩定市面，恢復社會秩序，具有重大之意義。

不應勒令停業

日人應以此為好

（中央社訊）不應勒令停業，日人應以此為好。據悉，政府正考慮對某些行業，勒令停業。然此舉，將對社會經濟，造成重大之影響。政府應慎重考慮，不應輕易勒令停業。同時，日人亦應以此為好，不應對中國之經濟，採取任何不利之行動。此項決定，已獲得政府之通過。

淪陷區公產

（中央社訊）淪陷區公產。據悉，政府正積極處理淪陷區之公產。此項公產，包括土地、房屋、工廠等。政府將根據國家法律，對這些公產，進行清理、變賣。所得之款項，將用於抗戰事業。此項決定，已獲得政府之通過。

陽床

（中央社訊）陽床。據悉，政府正積極推廣陽床。此種陽床，具有優良之性能，且價格低廉。政府將採取各種措施，鼓勵民眾購買陽床。此項決定，已獲得政府之通過。

Project introduction: <https://tinyurl.com/ecpo-intro>

Towards full text - 1

First steps

Expanding data: towards fulltext

- Manual typing not feasible
- Professional double-keying very expensive
- Challenges for OCR
 - Image: secondary copies - noisy, stains, scratches, etc
 - Document: dense layout, normal segmentation fails
 - Characters: special characters (emphasis), handwriting



[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](#)

世惡夫捕風捉影附會無
意不經之言名惟徒亂人
責者代為分謗甯非一
也哉戲填小詞以懲此報

of use

Source

還捏紛紛電報
塗說文章東抹
是一團糟墨白
擲壁最能虛造
市之大戲填小
意且令負責者
根之事張自不
焉所用惡惟惡
良神社會
惡消聽也
不辭勞
息過聰
非亦是
往此報
影附會無
符^^察一溢眾
通信社為報賊
排胎並字魏總也

排胎並字魏總也
通信社為報賊補助
符^^察一溢眾姑竝敵之不暇

Report bad result

俳辭一々（並序）惡術聽也

通信社為報腋秭助餞關良裨社鈺

苟々々察一泓庵々々郊謗之不暇

忍所用惡惟惡夫捕風從影附會無

根之事張皇不經之言名份徒亂人

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

滄說文帘束抹。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

滄說文章束抹 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

意且令

市之大

云爾

俳

通信社

苟能信

焉所用

根之事

意且令

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

滄說文章束抹 鈔郵花浪貼不辭勞

還捏紛紛電報

ca. 63% correctly recognized out-of-the-box (Abby FineReader)

Towards full text - 2

Segmentation

報晶

The Crystal

發行人 A.L. Teodoro

八二四四八號 號九二二號海山街

號三五九三〇九路電

家專影攝流一第上海
館相照術藝泰國
化裝貴相照化濟經價定
機良影攝光春此際
號五五路京州址地
號三九三〇九路電



專辦 高尙粵茶 大小適中 隨意小酌
華貴禮堂 租費從廉

杏花酒樓

路馬四 口里前豐
九四六二九

日三月三年八國民於新華本 厘四分貳幣國幣張貳日今 五期星 日一廿月四年捌拾貳國民華中

各地新聞 國內大事 國際形勢 社會生活 體育消息 藝術動態 科學新知 兒童樂園 婦女世界 兒童樂園 婦女世界



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
綏首重要據點多處已收復

（中央社訊）贛南各路華軍，自三月下旬起，大舉反攻，克復石龍、包圍南昌。高安方面，發生激烈爭奪戰，我軍奮勇殺敵，現正向高安推進中。綏首重要據點，多處已收復，我軍正積極清剿殘敵中。

抱定抗戰決心

中國現決不談和

英當局亦聲明未提議調停

（中央社訊）中國政府抱定抗戰決心，現決不談和。英當局亦聲明，未提議調停。中國政府表示，抗戰必能取得最後勝利，絕不考慮任何條件之停戰。

本報復刊小啓

本報自三月廿一日停刊後，承蒙各界人士厚愛，復刊後更將力求進步，務使內容充實，報導詳盡，以期服務社會，裨益人心。此啓。

陽痿早洩

可以早日治愈
保能如期

本藥專治陽痿早洩，功效顯著，服後立見奇效。各大藥房均有代售。

陽痿早洩

可以早日治愈
保能如期

本藥專治陽痿早洩，功效顯著，服後立見奇效。各大藥房均有代售。

各商店運動 今日可復業

被捕市民昨日均已釋放

（中央社訊）上海各商店運動，今日可復業。被捕市民昨日均已釋放。各界人士對政府之處理表示滿意。

河內 航船已過

（中央社訊）河內航船，昨日已過。航運恢復，對經濟發展有積極影響。

外交界動向

（中央社訊）外交界動向，近日有變。國際形勢緊張，各國正密切關注事態發展。

應有懸國旗之自由

（中央社訊）應有懸國旗之自由。公民應享有此項基本權利，政府應予以保障。

淪陷區公產

（中央社訊）淪陷區公產，政府正積極處理。旨在保護國家財產，維護社會穩定。

櫻餅

除皮清涼 平頂止痛 治頭痛 治牙痛 治胃痛 治腹痛

順風牌

鮮桔水 汽水 汽水 汽水

品出司公水汽華美

本報志實報道新聞 祝人人要說的話文字字淺近趣味濃厚是最平民化的刊物

本報具有廿一年的歷史是小型報紙其組織路普通全國刊登廣告最有效力

本報具有廿一年的歷史是小型報紙其粗銷路普遍全國刊登廣告最有力

法租界惡棍糾紛

今日可復業

被捕市民昨日均已釋放

文涉經過

法租界惡棍糾紛，經各方調解，今日可復業。被逮捕之市民，昨日均已釋放。此項糾紛，係因法租界內之惡棍，糾集多人，在公共場所滋事，擾亂治安。經政府派員調解，雙方達成協議，今日即可復業。被逮捕之市民，昨日均已釋放，現正接受調查。

河內

河內消息：法租界惡棍糾紛，今日可復業。被逮捕之市民，昨日均已釋放。此項糾紛，係因法租界內之惡棍，糾集多人，在公共場所滋事，擾亂治安。經政府派員調解，雙方達成協議，今日即可復業。被逮捕之市民，昨日均已釋放，現正接受調查。

已完

應有應國旗之自由

應有應國旗之自由，此為我國國民之基本權利。在淪陷區，國民應如何行使此項權利，實為當前之急務。本報特就此問題，邀請專家學者，進行深入探討。希望各界人士，能踴躍投稿，共同為維護國家尊嚴與國民自由而努力。

不堪勒索停業

不堪勒索停業，此為淪陷區商民之普遍遭遇。日寇在淪陷區實行高壓統治，對商民進行無休止的勒索，致使許多商號不堪重負，被迫停業。此種現象，嚴重損害了淪陷區之經濟發展，也反映了日寇殘暴統治之真面目。

一重天

一重天，淪陷區之黑暗與絕望。在日寇的鐵蹄下，淪陷區的天空陰沉無比，充滿了恐懼與絕望。商民生活在水深火熱之中，苦不堪言。此種黑暗之統治，必將導致淪陷區之徹底淪亡。

淪陷區公產
淪陷區公產，係指在淪陷區內之公有財產。其管理與處置，應符合國家利益與淪陷區之實際需要。本報將關注此項問題，並為公眾提供相關資訊。

杏花酒樓
高貴粵菜
華貴禮堂
租費從廉
大小適府
隨意小酌
九四六二九

好美華學園
高貴粵菜
華貴禮堂
租費從廉
大小適府
隨意小酌
九四六二九

陽痿早洩
可以早洩
治限日愈
治愈早洩
耐久保能

松蘭
治限日愈
治愈早洩
耐久保能

順風牌
汽水
鮮桔水
汽水
汽水

品需必之期暑
順風牌
汽水
鮮桔水
汽水
汽水

報晶
The Crystal
發行人 A. L. Teodoro
八二四三六四號 號八九二路日萬延里
新華之送寄單和地址詳印本報每份均送報中
《地址詳印本報每份均送報中》
《號五十九百七千三第》

家專影攝流一第
館相照術藝泰國
化裝貴相照化濟經預定
機良影攝光春此際
號五五路京州址地
號三九三〇九話電

華貴禮堂
租費從廉
大小適府
隨意小酌
九四六二九

日三月三年八國民華中 版四分第報報張日今 五期星 日一廿月四年八十二國民華中

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
緩管重要據點多處已收復



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，緩管重要據點多處已收復。此項軍事進展，顯示了我軍之強大戰鬥力，也反映了日寇在淪陷區之統治已趨於崩潰。我軍將繼續擴大戰果，爭取早日收復失地。

抱定抗戰決心
中國現決不談和

抱定抗戰決心，中國現決不談和。在淪陷區，我同胞應如何堅持抗戰，實為當前之急務。本報特就此問題，邀請專家學者，進行深入探討。希望各界人士，能踴躍投稿，共同為維護國家尊嚴與國民自由而努力。

本報復刊小啓

本報復刊小啓，本報自創刊以來，承蒙各界人士之愛護與支持，得以不斷發展。現因故停刊一段時間，現已籌備就緒，定於近日復刊。特此公告，敬請繼續關注。

本報志實報進新聞說人人要說的結文字淺近趣味濃厚是最平民化的刊物



- Home
- Schedule
- Organisers
- Register
- Resources / Details
- Submissions
- Results

RDCL2019 ICDAR Competition on Recognition of Documents with Complex Layouts

Overview

The competition presents challenges for page segmentation, region classification, and text recognition in an end-to-end scenario. The dataset contains scanned pages from contemporary magazines and technical articles. Participants will be provided with know-how and tools that aid the development or extension of their page analysis systems.

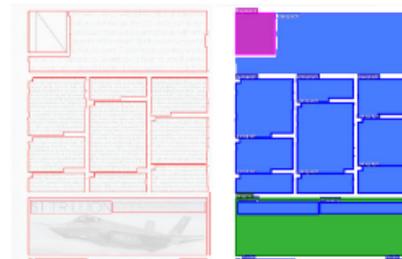
Dataset

Pages of magazines and technical articles



Main challenge

Page segmentation and region classification



<https://www.primaresearch.org/RDCL2019/>

**Document layout recognition and
page segmentation:**

**Page structuring elements -
separators**

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

《號四廿字〇報記登部工報和共公》

《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，自三月廿日起，大舉反攻，克復石龍、包圍南昌。高安發生激烈爭奪戰，緩督重要據點多處已收復。據悉：贛南各路華軍，自三月廿日起，大舉反攻，克復石龍、包圍南昌。高安發生激烈爭奪戰，緩督重要據點多處已收復。據悉：贛南各路華軍，自三月廿日起，大舉反攻，克復石龍、包圍南昌。高安發生激烈爭奪戰，緩督重要據點多處已收復。

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。英當局亦聲明未提議調停。據悉：中國政府抱定抗戰決心，現決不談和。英當局亦聲明未提議調停。據悉：中國政府抱定抗戰決心，現決不談和。英當局亦聲明未提議調停。

宿潼流寇各得勝

華軍甚得手

開慶軍戰北蘇

【本報廿日專電】宿潼流寇各得勝，華軍甚得手。開慶軍戰北蘇。據悉：宿潼流寇各得勝，華軍甚得手。開慶軍戰北蘇。據悉：宿潼流寇各得勝，華軍甚得手。開慶軍戰北蘇。

德機襲津

津市受驚

德機襲津

【本報廿日專電】德機襲津，津市受驚。據悉：德機襲津，津市受驚。據悉：德機襲津，津市受驚。

德機襲津

津市受驚

德機襲津

【本報廿日專電】德機襲津，津市受驚。據悉：德機襲津，津市受驚。據悉：德機襲津，津市受驚。

梅吳揚揚

梅吳揚揚

梅吳揚揚

【本報廿日專電】梅吳揚揚。據悉：梅吳揚揚。

杏花酒樓

華貴禮堂 租費從廉

高尙粵菜 大小筵席 隨意小酌

路馬四 里師德

九四六二九

高尙野味

妙手調出 色香味俱全

娛樂 妙手調出 色香味俱全

陽痿

治癒可以早洩

早洩 治癒可以早洩

本報復刊小啓

本報自一月一日起，由美商發行，恢復舊觀。今後本報有二大新向：促進東西南北民主團結，實現全全，以求民濟思益。實現此大之舉，取公開態度，為大輿論，凡有專事、小品、新聞、而於解脫民生痛苦，進入羣眾。尤敢意，此其一二。請希留意。

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

【本報廿日專電】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。

河內 航郵已通

【本報廿日專電】河內航郵已通。據悉：河內航郵已通。

外交界動靜

【本報廿日專電】外交界動靜。據悉：外交界動靜。

不推勒索停業

星屬島震心病狂

與日人服比為奸

【本報廿日專電】不推勒索停業。星屬島震心病狂。與日人服比為奸。

任何國民 應有懸國旗之自由

【本報廿日專電】任何國民，應有懸國旗之自由。據悉：任何國民，應有懸國旗之自由。

淪陷區公產

【本報廿日專電】淪陷區公產。據悉：淪陷區公產。

小說「重天」

【本報廿日專電】小說「重天」。據悉：小說「重天」。

淪陷區公產

【本報廿日專電】淪陷區公產。據悉：淪陷區公產。

紛糾旗懸界租法

【本報廿日專電】紛糾旗懸界租法。據悉：紛糾旗懸界租法。

交通經過

【本報廿日專電】交通經過。據悉：交通經過。

河內 航郵已通

【本報廿日專電】河內航郵已通。據悉：河內航郵已通。

不推勒索停業

星屬島震心病狂

與日人服比為奸

【本報廿日專電】不推勒索停業。星屬島震心病狂。與日人服比為奸。

任何國民 應有懸國旗之自由

【本報廿日專電】任何國民，應有懸國旗之自由。據悉：任何國民，應有懸國旗之自由。

淪陷區公產

【本報廿日專電】淪陷區公產。據悉：淪陷區公產。

小說「重天」

【本報廿日專電】小說「重天」。據悉：小說「重天」。

淪陷區公產

【本報廿日專電】淪陷區公產。據悉：淪陷區公產。

本報具有廿一年的歷史是小型報紙鼻祖銷路普遍全國刊登廣告最有力

暑之期必之需品

順風牌 汽水 汽水 汽水

四大特點：鮮美、清潔、衛生、可口

上海華華汽水公司出品

**Document layout recognition and
page segmentation:**

**Page structuring elements –
registers**

報晶

The Crystal

發行人 A. L. Teodoro

八二四三九路電 號九九二路口漢址館
紙帶之送寄益便包總照於登立號街進特政郵華中

《號四廿字 〇證登報部工具和共公》

《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
號五五路京南址地
號三九三〇九話電



路馬四口里領豐
辦專
高尙粵茶
華貴禮堂
租費從廉

杏花酒樓

大小筵席
隨意酌酌
隨小酌
隨大筵席
九四六二九

日三月三年八國民於前報本 厘四分貳幣國售張貳日今 五期星 日一廿月四年捌拾貳國民華中

1 reg - 8 char

（本報訊）...



4 reg - 35 char

5 reg - (44 char)

克復石叻包圍南昌

0.5 reg - 4 char

抱定抗戰決心
中國現決不談和

3 reg - 26 char

華軍甚得手

1 reg - 8 char

本報復刊小啓

1 reg - 8 char

本報復刊小啓

各路華軍大反攻
克復石叻包圍南昌
（本報訊）...

（本報訊）...

（本報訊）...

（本報訊）...

（本報訊）...

（本報訊）...

園今美好

高尙妙演劇團
換樂妙演劇團
園地之演樂

陽痿
治限可以
愈日早洩
如保
耐如

小便風濕... 遺精... 淋病...

櫻敵
除痰止咳
平喘止咳
補肺防癆
益肺防癆

品出廠膠保記信泰正海上海

品出司公水汽華美

品需必之期暑

點特大四
清利鮮潤
潔學美清
衛製可蒸
生造口溫

汽鮮順
桔風
水水牌

品出司公水汽華美

今日可復業

被捕市民昨日均已釋放

河內 航郵已通

1 reg - 8 char

2 reg - 17 char

3 reg - 26 char

1.5 reg - 12 char

3 reg - (26 char)

1.5 reg - 12 char

1 reg - 8 char

2 reg - 17 char

4 reg - 35 char

5 reg - (44 char)

0.5 reg - 4 char

3 reg - 26 char

1 reg - 8 char

1 reg - 8 char

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

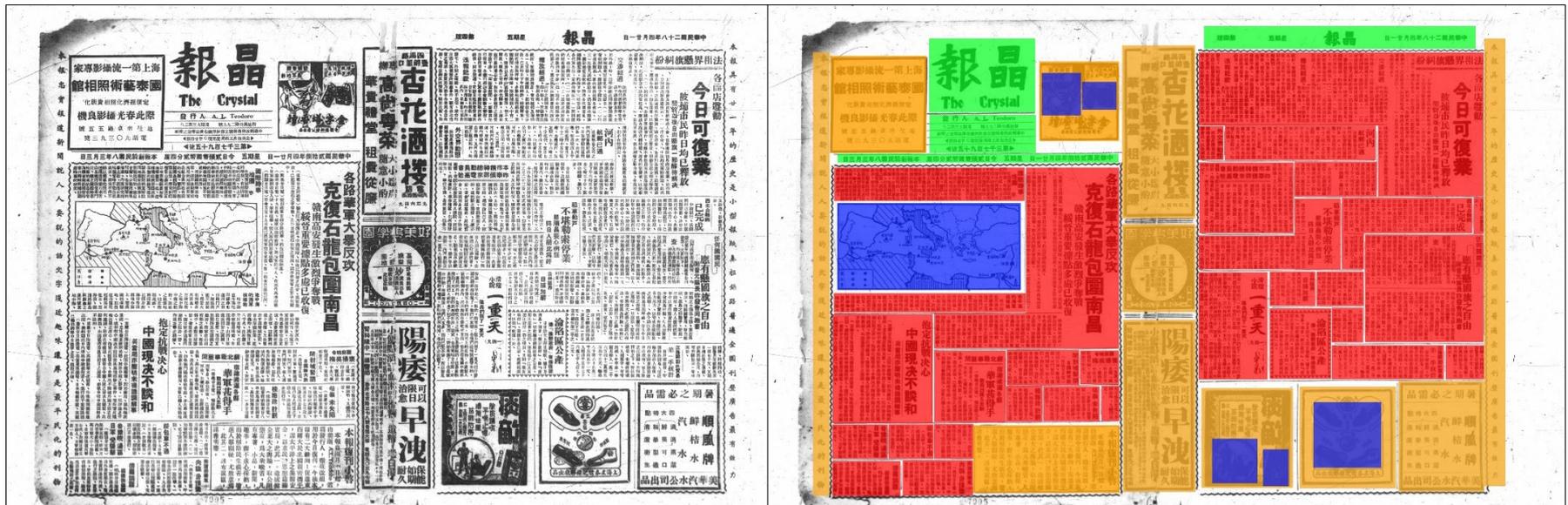
本報具有廿一年的歷史是小型報紙鼻祖銷路普遍全國刊登廣告最有力

**Document layout recognition and
page segmentation:**

Crowdsourcing

Crowdsourcing segmentation

- Pilot project with Pallas Ludens GmbH
- Let the (controlled) crowd help analyzing the pages
 - Identify and label four item types:
 - image/drawing
 - article
 - advertisement
 - additional information
 - Closely monitored and supervised
 - Crowd members NOT able to speak/read Chinese



晶報

第... 廣告... 例刊...

定價... 零售... 廣告...

科發白濁丸

專治白濁... 功效神速...



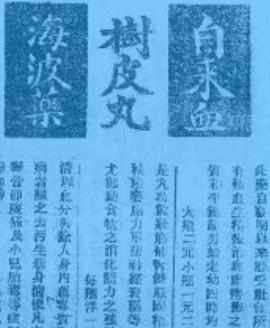
請吸新出 煙香貨國牌城長大

南洋兄弟煙草公司



上海五洲大藥房

自來血 樹皮丸 海波藥



中法儲蓄會

儲蓄世界 金銀如流 水東來西 去欲省而 不能今有 妙法每月 在消耗費 中節省一 份錢去投 入儲蓄會 裏將來可 以集成巨 數并且每 月有得獎 的希望中 法儲蓄會 成立最早 會員最多 營業開張 最廣得任 中國政府 立案注冊 政府注冊 開獎在即 不要錯過

德六零六

專治梅毒... 功效顯著...



西醫張世楷

專治... 醫術精湛...

美少年

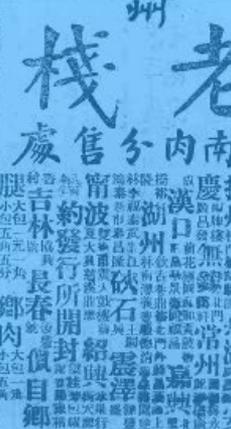
美容... 恢復青春...

秋令咳嗽

止咳金丹 天賦靈藥 製此靈藥 三分鐘內 咳止嗽停 貧寒而索 不取分文 富貴人本 祇收藥銀 每瓶一元 不靈還銀 上海英租界... 造福齋藥房

預防花柳奇藥 愛克憐

德國六零六 只能醫梅毒 新出愛克憐 真是藥中神 既能防梅毒 又可防五淋 諸君如用此 消患於無形



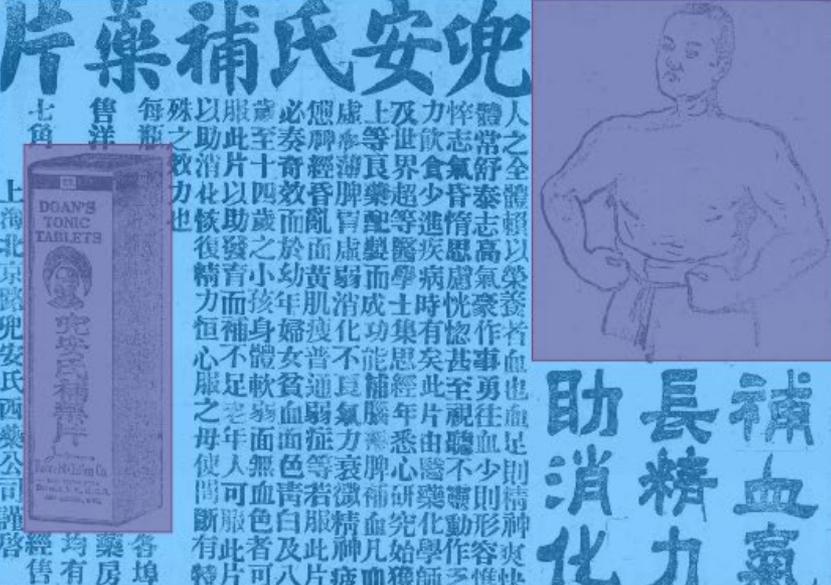
萬隆老棧

冬令特設 馳名肉分售處

南京 蘇州 揚州 無錫 常州 鎮江 南通 杭州 寧波 紹興 嘉興 湖州 溫州 台州 處州 衢州 嚴州 金華 衢州 嚴州 金華

兜安氏補藥片

補血氣 助消化 長精力

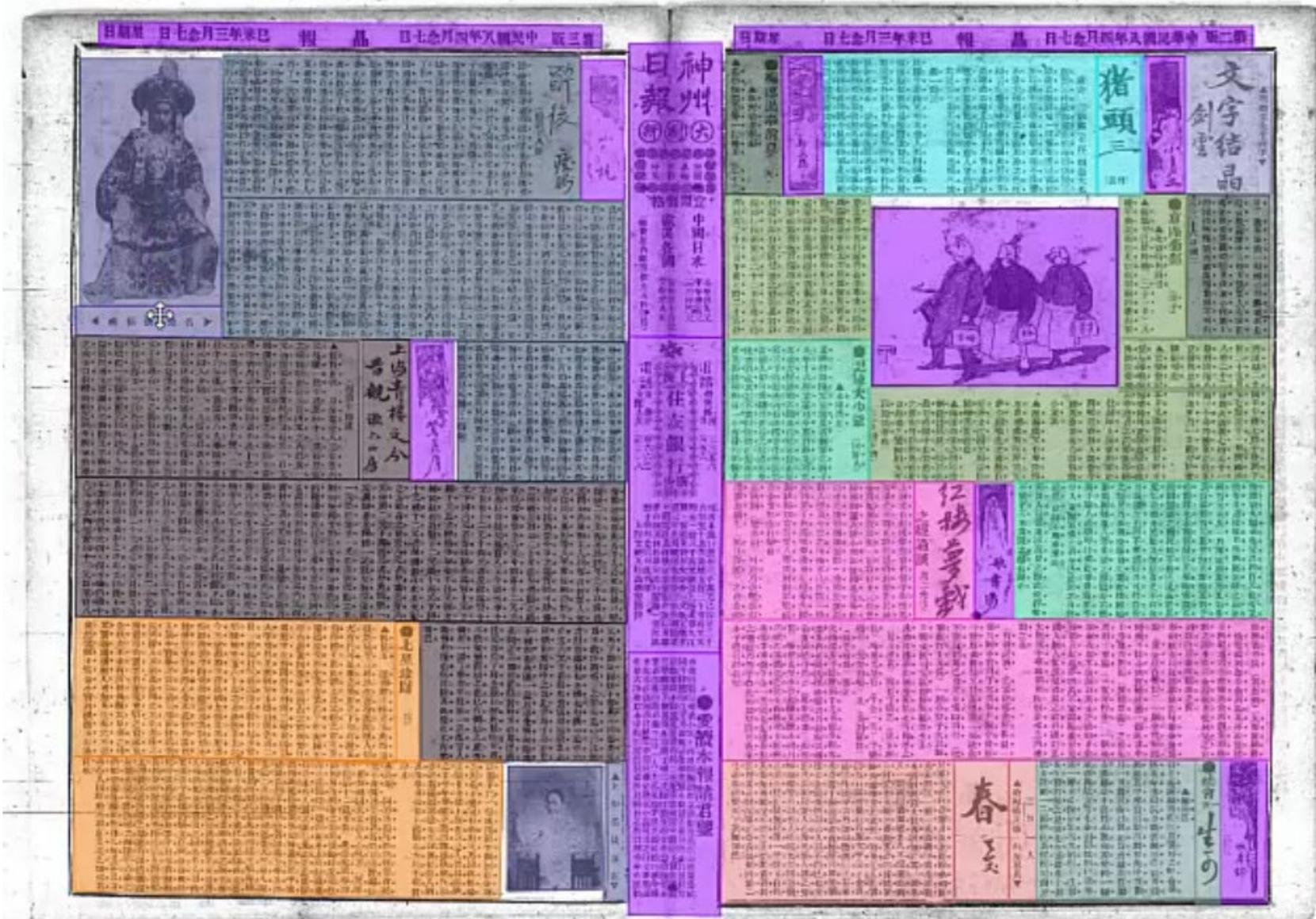


DOAN'S TONIC TABLETS

寒雲主人好古知書深得三代漢魏之神髓... 求書者不暇爰擬定書例如下...

中華郵政特准掛號認爲新聞紙類... 中華郵政特准掛號立券之報紙... 已在中華郵政特准掛號認爲新聞紙類...

Grouping semantic units



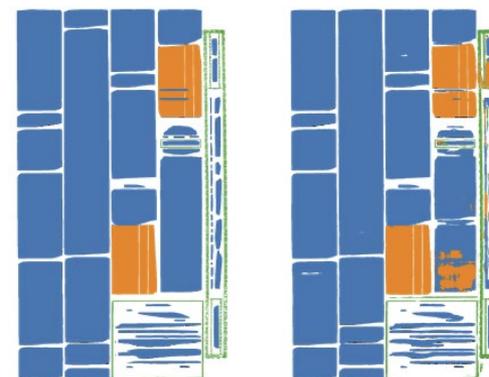
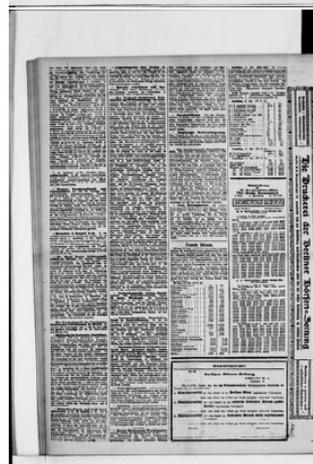
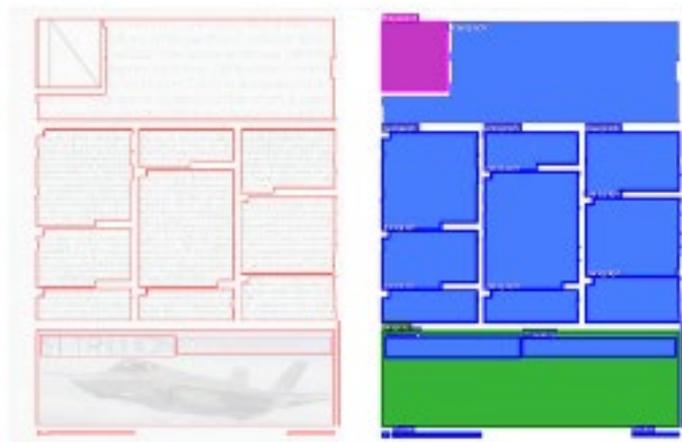
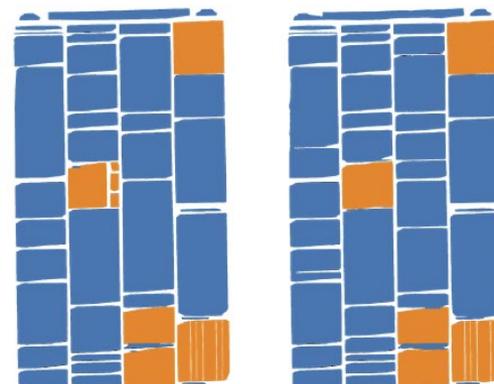
Outcome of segmentation pilot

1. Page segmentation can be outsourced to expert crowd
 - Requires supervision
 - Advanced user interfaces (high usability, efficiency)
 - For more advanced tasks: Crowd should read Chinese (semantic grouping)
2. *Jingbao* 晶報 1919-21 (927 folds) completely segmented with qualified boxes, issues of April 1919 with semantic units
3. Further processing:
 - Computational Knowledge Lab (知識計算實驗室), Dept. Engineering Science and Ocean Engineering, Taiwan National University, <http://www.cklab.org/>
 - Advanced Linguistic Technologies, Heidelberg
 - Seeking partners for collaboration!

Towards full text - 3:

Machine learning

Neural networks



ICDAR 2019 – RDCL2019 challenge

Liebl and Burghardt 2020

Neural networks

Large research projects on (western, mostly Latin script) newspapers, three examples from Europe:

- Europeana newspapers <http://www.europeana-newspapers.eu/>
- Impresso <https://impresso-project.ch/>
- NewsEye <https://www.newseye.eu/>
- OCR-D <https://ocr-d.de/>

Re-usable (free, open source, open access) software, for example dhSegment - a “generic deep-learning framework for Historical Document Processing” <https://github.com/dhlab-epfl/dhSegment>

Our aims:

- Use for segmentation
- Train models to detect main element classes:
Paratext (headers and marginalia), Advertisement and images, Text areas
- Requirement: Ground Truth

Towards full text - 4:

Ground truth

Funding for Ground Truth

Short term funding by Field of Focus 3 (ExStrat Heidelberg)

Small team with Duncan Paterson and 4 RA's (typing)

- GroundTruth 1: Bounding boxes with semantic labels
 - Coordinates in json and web-annotation format
 - Jingbao April 1939 und April 1938
 - Outcome: 70 folds, 6335 shapes, Ø 90,5 shapes/fold

Annotation tool



mode.select

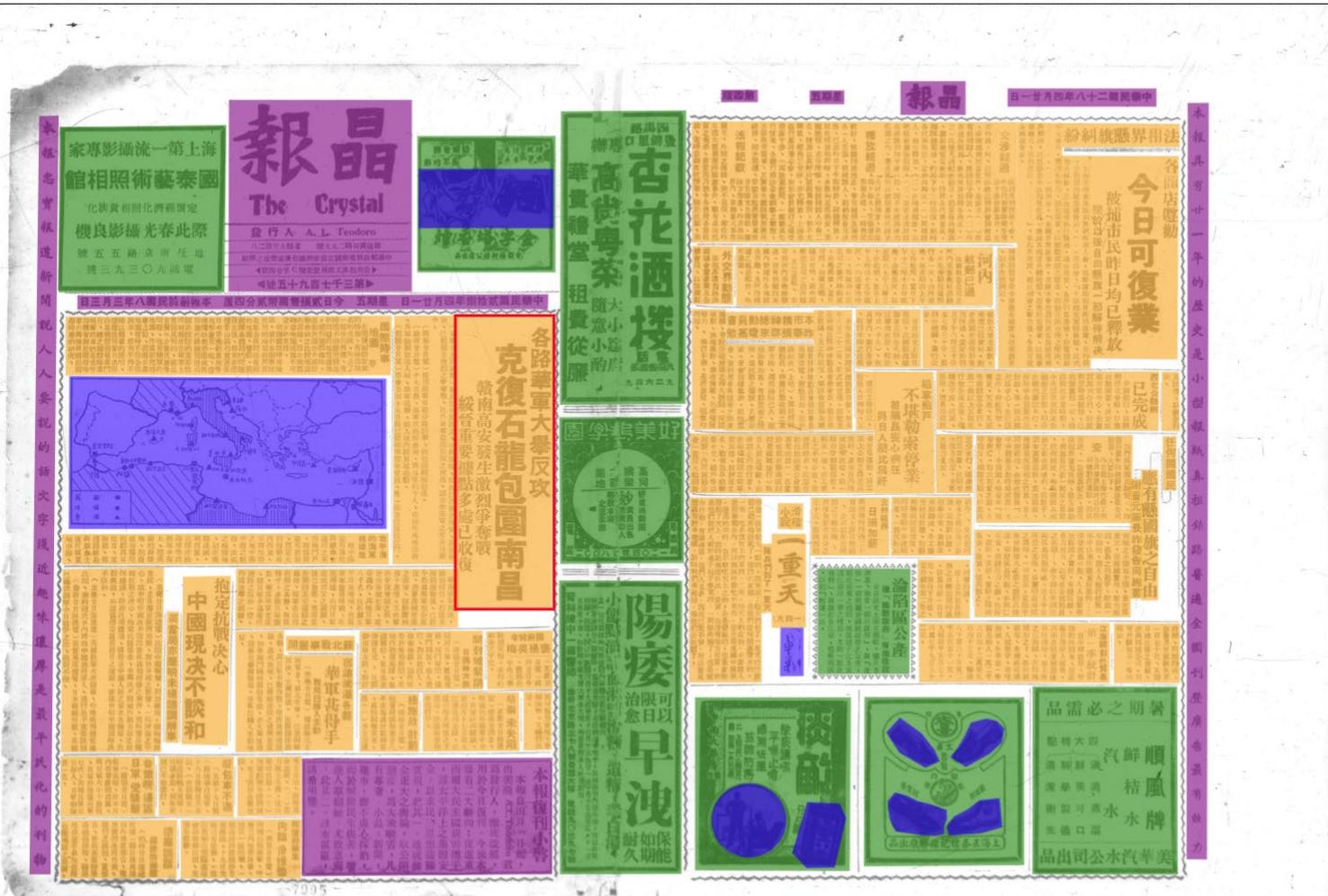


Image or Drawing

Article

Advertisement

Additional Information

SUBMIT



Annotation output 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 }
```

Funding for Ground Truth - II

Short term funding by Field of Focus 3 (ExStrat Heidelberg)

Small team with Duncan Paterson and 4 RA's (typing)

- GroundTruth 1: Bounding boxes with semantic labels
 - Coordinates in json and web-annotation format
 - Jingbao April 1939 und April 1938
 - Outcome: 70 folds, 6335 shapes, Ø 90,5 shapes/fold
- GroundTruth 2: Texts
 - Blind double-keying
 - Jingbao April 1939 - 10 issues, 40 folds
 - Outcome: ca. 245.000 characters, local XML format, Ø 6100 chars/fold (typing slow ~10h/fold)

GT full text

Local XML schema

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

Encoding in “UTF-8”:

- Characters according to original, wherever the respective Unicode code point was available
- Illegible: `&gaiji;`
- All characters entered in double-space form

Special feature in East Asian texts: different running directions

- `<div>` with `@mode` and `@dir`
- Default: `mode="vertical-rl"`

<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>
```

```

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>

```

Default setting per page:
 @mode=„vertical-rl“ and
 @dir=„rtl“ (right-to-left)

```

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>

```

If needed, we encode changes in sub-level <div>'s
 After </div> all changes reset to default

本報具有廿一年的歷史是小型報紙鼻祖銷路普遍全國刊登廣告最有力

各商店應勸

今日可復業

被捕市民昨日均已釋放

關於以前被捕市民，經各界再三交涉，業經有關機關，於昨日全部釋放，並發給津貼，以資生活。此項被捕市民，係在日前，因參加抗戰，而被捕者。現已獲釋，並發給津貼，以資生活。此項被捕市民，係在日前，因參加抗戰，而被捕者。現已獲釋，並發給津貼，以資生活。

紛糾旗懸界用法

交涉難題

關於懸旗問題，各界人士，均表反對。認為懸旗，係屬非法行為，應予禁止。且懸旗，易引起國際糾紛，對國家利益，殊屬不利。故應嚴禁懸旗，以維國家尊嚴。此項懸旗問題，現正由有關機關，進行交涉中。各界人士，應予諒解，並配合有關機關，共同維護國家利益。

河內

河內方面，近日局勢緊張。各方軍隊，紛紛調動，以資防範。且河內，為交通要道，各方勢力，均欲爭奪。故河內局勢，極為複雜。各方人士，應密切注意河內局勢之發展，並根據實際情況，採取適當之行動。

外交界動靜

外交界近日動靜頻繁。各方代表，紛紛抵滬，進行會談。且外交界，正就各項問題，進行磋商。各方人士，應密切注意外交界之動靜，並根據實際情況，採取適當之行動。

不推勒索停業

近日勒索事件，層出不窮。各方人士，均表不滿。認為勒索，係屬非法行為，應予嚴懲。且勒索，易引起社會動亂，對國家利益，殊屬不利。故應嚴禁勒索，以維社會安寧。此項勒索事件，現正由有關機關，進行調查中。各界人士，應予配合，共同維護社會安寧。

應有懸國旗之自由

應有懸國旗之自由，為國民之基本權利。政府應予保障，不得任意剝奪。且懸國旗，為愛國之表現，應予鼓勵。故政府應採取措施，保障國民懸國旗之自由。此項懸國旗之自由，現正由有關機關，進行研究。各界人士，應予支持，共同維護國民之基本權利。

滬陷區公產

滬陷區公產，現正由有關機關，進行清理。各方人士，應予配合，共同維護公產之安全。且公產，為國家之財產，應予嚴加保護。故有關機關，應採取措施，防止公產之流失。此項公產清理工作，現正由有關機關，進行中。各界人士，應予支持，共同維護公產之安全。

一重天

一重天，為抗戰之精神。各方人士，應以此為楷模，奮發圖強。且一重天，為抗戰之動力，應予弘揚。故有關機關，應採取措施，弘揚一重天精神。此項一重天精神，現正由有關機關，進行推廣。各界人士，應予支持，共同弘揚一重天精神。

陽痿早洩

陽痿早洩，為常見之病症。各方人士，應予重視，並採取適當之治療。且陽痿早洩，易引起生活痛苦，應予預防。故有關機關，應採取措施，預防陽痿早洩。此項陽痿早洩之治療，現正由有關機關，進行研究。各界人士，應予支持，共同預防陽痿早洩。

克復石龍包圍南昌

克復石龍，包圍南昌，為抗戰之重要戰役。各方人士，應以此為契機，奮發圖強。且石龍，為抗戰之要地，應予嚴加保護。故有關機關，應採取措施，克復石龍。此項石龍克復工作，現正由有關機關，進行中。各界人士，應予支持，共同克復石龍。

各路華軍大舉反攻

各路華軍，大舉反攻，為抗戰之重要戰役。各方人士，應以此為契機，奮發圖強。且華軍，為抗戰之主力，應予嚴加保護。故有關機關，應採取措施，大舉反攻。此項華軍反攻工作，現正由有關機關，進行中。各界人士，應予支持，共同大舉反攻。

贛南高安發生激烈爭奪戰

贛南高安，發生激烈爭奪戰。各方人士，應以此為契機，奮發圖強。且高安，為抗戰之要地，應予嚴加保護。故有關機關，應採取措施，克復高安。此項高安克復工作，現正由有關機關，進行中。各界人士，應予支持，共同克復高安。

抱定抗戰決心

抱定抗戰決心，為抗戰之精神。各方人士，應以此為楷模，奮發圖強。且抗戰，為民族之存亡，應予嚴加保護。故有關機關，應採取措施，抱定抗戰決心。此項抗戰決心，現正由有關機關，進行推廣。各界人士，應予支持，共同抱定抗戰決心。

中國現決不談和

中國現決不談和，為抗戰之決心。各方人士，應以此為楷模，奮發圖強。且談和，為民族之恥，應予嚴加保護。故有關機關，應採取措施，決不談和。此項決不談和之決心，現正由有關機關，進行推廣。各界人士，應予支持，共同決不談和。

前線華軍北進

前線華軍，北進抗戰。各方人士，應以此為契機，奮發圖強。且華軍，為抗戰之主力，應予嚴加保護。故有關機關，應採取措施，北進抗戰。此項華軍北進工作，現正由有關機關，進行中。各界人士，應予支持，共同北進抗戰。

本報復刊小啟

本報復刊小啟，為本報之重要公告。各方人士，應以此為契機，奮發圖強。且本報，為抗戰之喉舌，應予嚴加保護。故有關機關，應採取措施，復刊本報。此項本報復刊工作，現正由有關機關，進行中。各界人士，應予支持，共同復刊本報。

單軍甚得手

單軍甚得手，為抗戰之重要戰役。各方人士，應以此為契機，奮發圖強。且單軍，為抗戰之主力，應予嚴加保護。故有關機關，應採取措施，單軍甚得手。此項單軍甚得手工作，現正由有關機關，進行中。各界人士，應予支持，共同單軍甚得手。

英軍常備軍將領

英軍常備軍將領，為抗戰之重要人物。各方人士，應以此為楷模，奮發圖強。且英軍，為抗戰之主力，應予嚴加保護。故有關機關，應採取措施，英軍常備軍將領。此項英軍常備軍將領工作，現正由有關機關，進行中。各界人士，應予支持，共同英軍常備軍將領。

本報復刊小啟

本報復刊小啟，為本報之重要公告。各方人士，應以此為契機，奮發圖強。且本報，為抗戰之喉舌，應予嚴加保護。故有關機關，應採取措施，復刊本報。此項本報復刊工作，現正由有關機關，進行中。各界人士，應予支持，共同復刊本報。

本報復刊小啟

本報復刊小啟，為本報之重要公告。各方人士，應以此為契機，奮發圖強。且本報，為抗戰之喉舌，應予嚴加保護。故有關機關，應採取措施，復刊本報。此項本報復刊工作，現正由有關機關，進行中。各界人士，應予支持，共同復刊本報。

本報復刊小啟

本報復刊小啟，為本報之重要公告。各方人士，應以此為契機，奮發圖強。且本報，為抗戰之喉舌，應予嚴加保護。故有關機關，應採取措施，復刊本報。此項本報復刊工作，現正由有關機關，進行中。各界人士，應予支持，共同復刊本報。

淡齋

治限可以
愈日早洩
耐如保能

陽痿早洩

治限可以
愈日早洩
耐如保能

本報復刊小啟

本報復刊小啟，為本報之重要公告。各方人士，應以此為契機，奮發圖強。且本報，為抗戰之喉舌，應予嚴加保護。故有關機關，應採取措施，復刊本報。此項本報復刊工作，現正由有關機關，進行中。各界人士，應予支持，共同復刊本報。

品需必之期暑

順風牌
鮮桔水
汽水
汽水
汽水

品需必之期暑

順風牌
鮮桔水
汽水
汽水
汽水

淡齋

治限可以
愈日早洩
耐如保能

家專影攝流一第上海

館相照術藝泰國

花樣實相照化濟經價定
機良影攝光春此際
舖三五路京市址地
舖三九三〇九話電

報晶

The Crystal

發行人 A. L. Teodoro
八二號三五路京市址地
舖三五路京市址地
舖三九三〇九話電

高尙粵菜

華貴禮堂
租費從廉

香花酒樓

高尙粵菜
租費從廉

克復石龍包圍南昌

各路華軍大舉反攻
贛南高安發生激烈爭奪戰
抱定抗戰決心
中國現決不談和

本報復刊小啟

本報復刊小啟，為本報之重要公告。各方人士，應以此為契機，奮發圖強。且本報，為抗戰之喉舌，應予嚴加保護。故有關機關，應採取措施，復刊本報。此項本報復刊工作，現正由有關機關，進行中。各界人士，應予支持，共同復刊本報。

本報忠實報導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

Towards full text - 5:

First results

報晶

The Crystal

發行人 A.L. Teodoro

八二四四八號 號九九二號海山街

號三五九三〇九路電

家專影攝流一第上海
館相照術藝泰國

化裝貴相照化濟經價定
機良影攝光春此際

號五五路京州址地
號三九三〇九路電



專辦 高尙粵茶 大小適中 隨意小酌 華貴禮堂 租費從廉

杏花酒樓

路馬四 口里師雙

九四六二九

日三月三年八國民於新華本 厘四分貳幣國幣張貳日今 五期星 日一廿月四年捌拾貳國民華中

克復石龍包圍南昌
贛南高安發生激烈爭奪戰
綏首重要據點多處已收復

（中央社訊）贛南高安發生激烈爭奪戰，我軍奮勇進攻，石龍、包圍、南昌等處重要據點多處已收復。敵軍傷亡慘重，正向南線潰退。我軍正乘勝追擊中。



（中央社訊）贛南高安發生激烈爭奪戰，我軍奮勇進攻，石龍、包圍、南昌等處重要據點多處已收復。敵軍傷亡慘重，正向南線潰退。我軍正乘勝追擊中。

抱定抗戰決心 中國現決不談和
英當局亦聲明未提議調解

（中央社訊）中國政府抱定抗戰決心，現決不談和。英當局亦聲明未提議調解。中國政府表示，抗戰必能取得最後勝利。

本報復刊小啓
本報自復刊以來，承蒙各界人士之厚愛，業務蒸蒸日上。現因業務需要，特將本報遷往新址辦公。特此啓事。

華軍甚得手
（中央社訊）華軍在戰場上表現英勇，甚得手。敵軍傷亡慘重，我軍正乘勝追擊中。

推展救濟計劃
（中央社訊）政府正積極推展救濟計劃，以幫助抗戰後之民生。各項救濟措施正積極落實中。

日本現決不談和
（中央社訊）日本現決不談和，其意圖何在，已昭然若揭。中國政府將繼續抗戰到底。

日本現決不談和
（中央社訊）日本現決不談和，其意圖何在，已昭然若揭。中國政府將繼續抗戰到底。

日本現決不談和
（中央社訊）日本現決不談和，其意圖何在，已昭然若揭。中國政府將繼續抗戰到底。

好美美
高尙粵茶 大小適中 隨意小酌 華貴禮堂 租費從廉

（中央社訊）好美美茶樓，環境優雅，服務周到。是親友聚會、商務洽談之理想場所。

陽痿
可以早日治愈
保能如期

（本報訊）陽痿之症，多由腎虛、氣弱所致。本藥房特製之補腎丸，功效顯著，能早日治愈。

早洩
保能如期

（本報訊）早洩之症，多由神經衰弱、精液不足所致。本藥房特製之補腎丸，功效顯著，能早日治愈。

法界懸旗糾紛 各商店邀勸 今日可復業

（中央社訊）法界懸旗糾紛，各商店邀勸，今日可復業。法界各界人士，為維護法界治安，特發起懸旗糾紛。各商店應予配合，共同維護法界秩序。

河內 航船已過

（中央社訊）河內航船已過，交通恢復。河內各界人士，對航船過河表示歡迎。交通恢復，對河內經濟發展將有重大貢獻。

不推勒索停業 吳日人以此為好

（中央社訊）不推勒索停業，吳日人以此為好。吳日人表示，不推勒索，是維護法界秩序之良舉。吳日人對此表示支持。

應有懸國旗之自由

（中央社訊）應有懸國旗之自由，是國民之基本權利。政府應保障國民之基本權利，維護國家尊嚴。

淪陷區公產

（中央社訊）淪陷區公產，應予保護。淪陷區公產，是國家財產之重要組成部分。政府應採取措施，保護淪陷區公產。

小報一重天

（中央社訊）小報一重天，是國民之精神食糧。小報應報導事實，傳播真理，為國民提供精神支持。

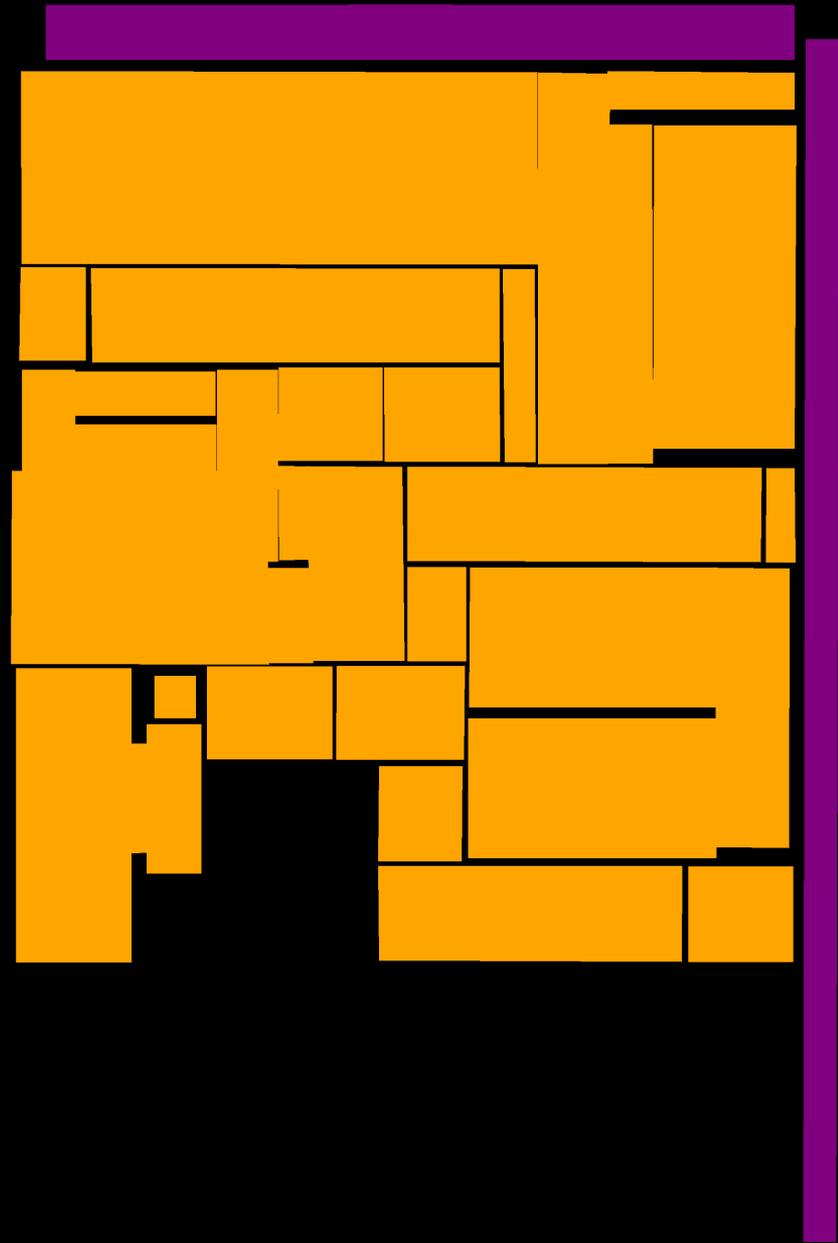
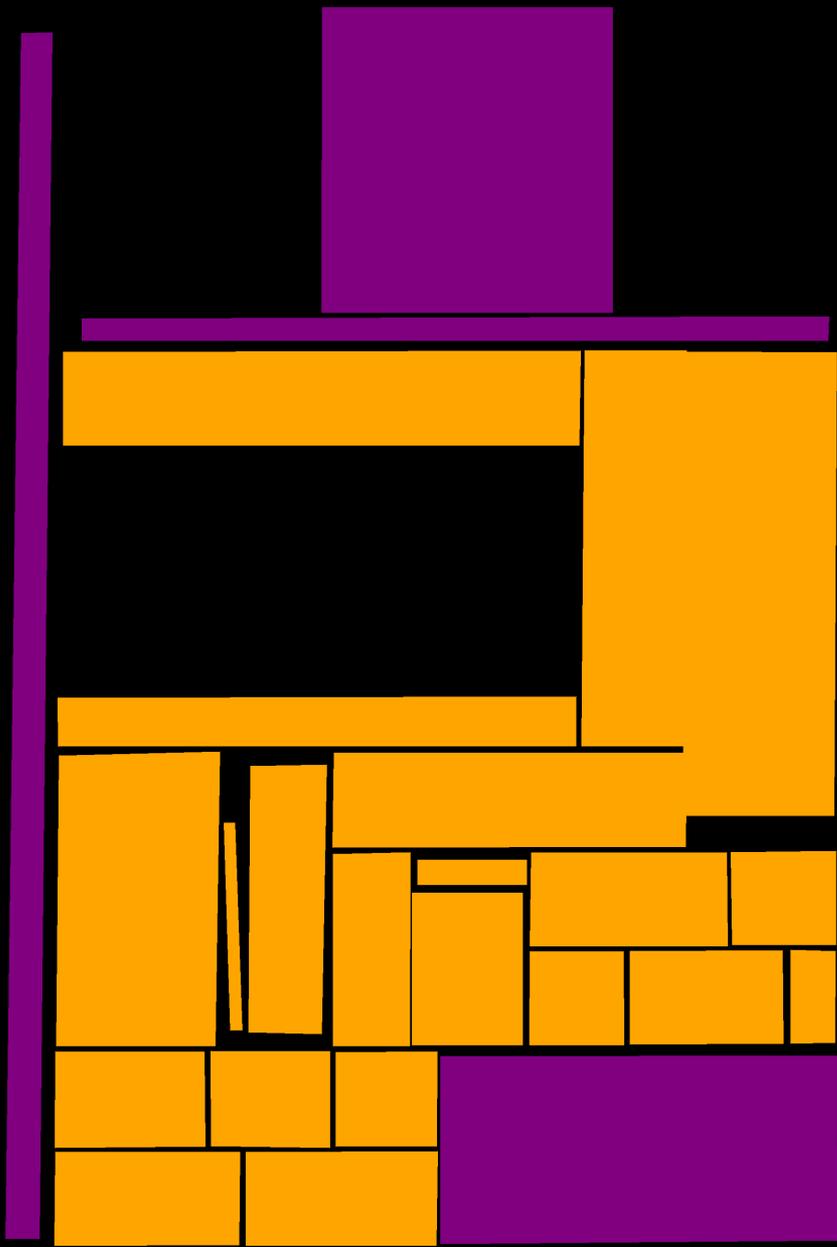
美華汽水公司出品

順風牌 鮮桔水 汽水 汽水 汽水

（本報訊）美華汽水公司出品之順風牌鮮桔水，口味清爽，營養豐富。是消暑解渴之佳品。

本報志實報道新聞 祝人人要說的話文字淺近趣味濃厚是最平民化的刊物

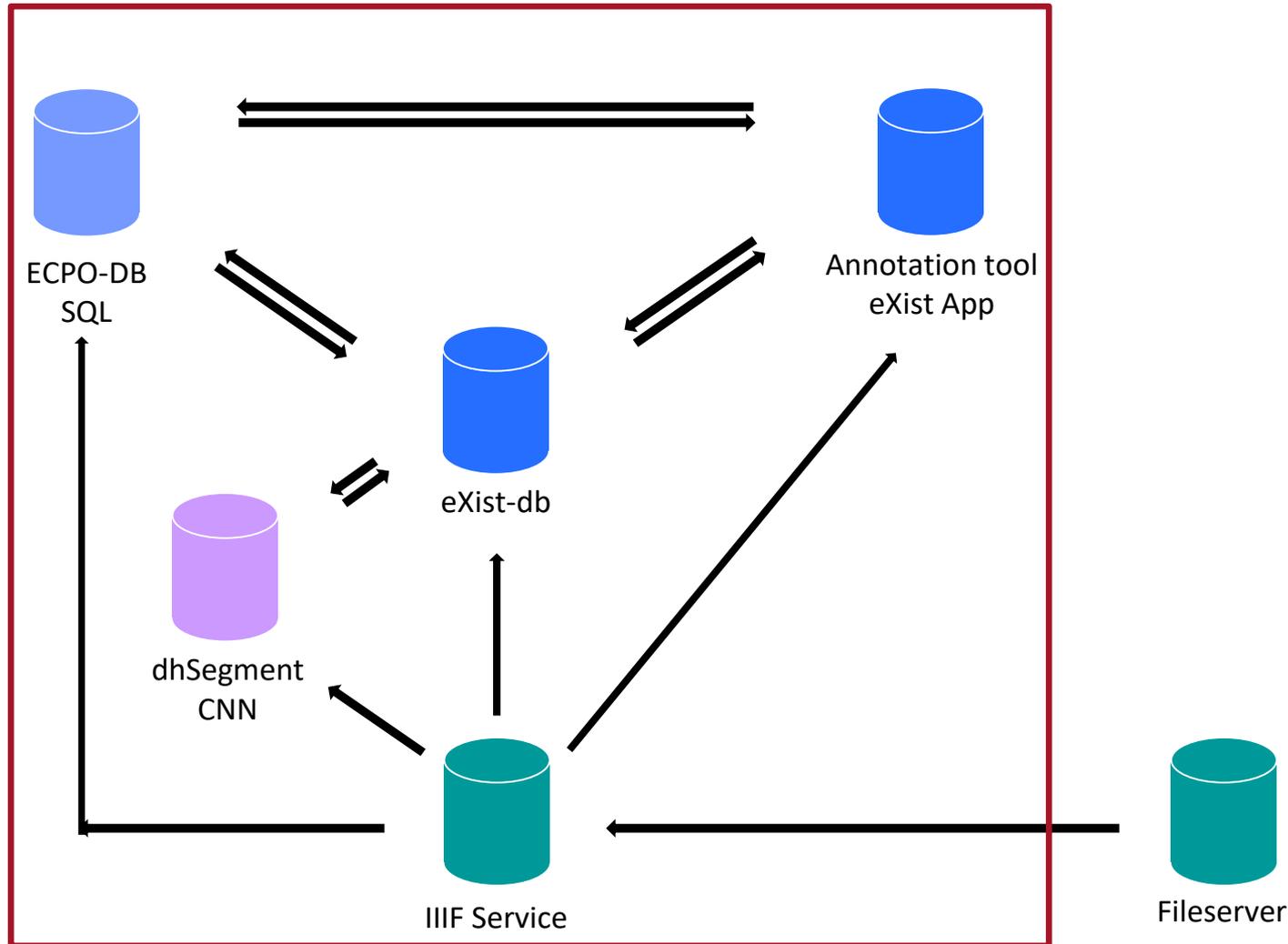
本報具有廿一年的歷史是小型報紙其組織路普通全國刊登廣告最有效力



Towards full text - 6:

Infrastructure

ECPO – Expanding the infrastructure



Closing remarks
(more challenges)

Challenges - I

Scalable infrastructure to accommodate additional modules, and (potentially) more data/traffic

- Solution: re-design of infrastructure, containerisation, XML-db (texts)

Neural networks require strong machines (e.g. GPU server)

- Solution: bwForCluster MLS&WISO (High Performance Cluster)

Looking for partners to collaborate

- Our data is available for re-use, if you feel challenged by Chinese newspapers from the Republican era – do go ahead
- If you are interested in any kind of partnership – do let me know

Challenges - II

Fighting the Latin-script bias

Algorithms working for European medieval manuscripts or 19th ct. Newspapers can NOT just be re-used

Especially non-Latin / non-alphabetic script material shows characteristics, existing software solutions may still be unaware of

Solutions?

publish data-sets, make these characteristics more commonly known, explain specific challenges to the expert-communities

Challenges - III

Fighting biased opinion - „Aren't the Chinese already doing it, anyway?“

- Little is known about Chinese research on processing historical material
- Research about Republican newspapers is no “hot topic“, like, for example, “handwritten text recognition”
- Academic projects using Neural networks for research are scarce
- Chinese Libraries usually “just” offer author-title indexes
- Access to academic research in China usually behind paywall and/or limited from outside China

=> Systematic literature overview (in progress)

- Where can relevant information be found
- Which projects use neural networks for their research
- Which approaches were chosen, what worked, what did not
- Can results be reproduced, are models and datasets shared
- Are we able to re-use their algorithms, models, data sets
- Our report will be submitted to a Special Issue on Digital Humanities and East Asian Studies, IJDH (due May 2021)

Our Team



Duncan Paterson



Jia Xie 谢佳



Matthias Arnold



Part of the larger research group

Selected 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.

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

Arnold, Matthias, and Lena Hessel. "Transforming Data Silos into Knowledge: Early Chinese Periodicals Online (ECPO)." In *Heuveline, Vincent, Gebhart, Fabian und Mohammadianbisheh, Nina (Hrsg.): E-Science-Tage 2019: Data to Knowledge*, 95–109. Heidelberg: heiBOOKS, 2020. <https://doi.org/10.11588/heibooks.598.c8420>.

Sung, Doris, Lying Sun, and 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 (2): 227–37, 2014. <https://doi.org/10.1353/tsw.2014.0004>.

Arnold, Matthias. "Multilingual Research Projects: Challenges for Making Use of Standards, Authority Files, and Character Recognition." *Digital Studies / Le champ numérique*, special issue "Towards Multilingualism in Digital Humanities", ed. Martin Lee und Cosima Wagner (manuscript accepted by Special issue editors, currently under review by journal).

Wick, Christoph, and Frank Puppe. "Fully Convolutional Neural Networks for Page Segmentation of Historical Document Images." *2018 13th IAPR International Workshop on Document Analysis Systems (DAS)*, April 2018, 287–92. doi:[10.1109/DAS.2018.39](https://doi.org/10.1109/DAS.2018.39)

Liebl, Bernhard, and Manuel Burghardt. "An Evaluation of DNN Architectures for Page Segmentation of Historical Newspapers," 15 April 2020. <https://arxiv.org/abs/2004.07317v1>

Early Chinese Periodicals Online (ECPO) – Project introduction. AAS2021. <https://tinyurl.com/ecpo-intro>

Contact

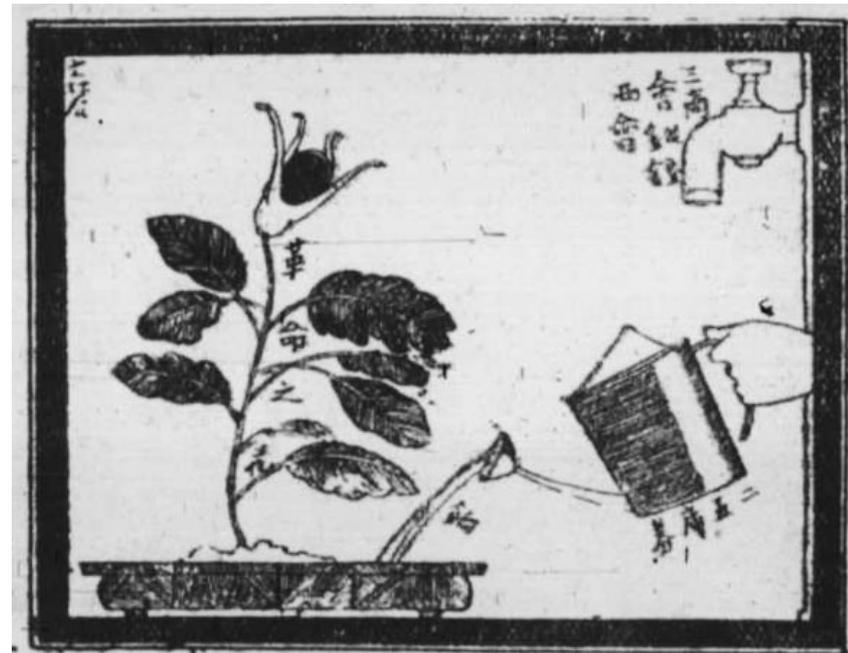
Matthias Arnold

matthias.arnold@uni-hd.de

<http://tinyurl.com/matthias-arnold>

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

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



Thank you!

