

Layout Parser, come creare un dataset di qualità per allenare l'Intelligenza Artificiale

Silvano Imboden, Gabriele Marconi, Simona Caraceni, Rossella Pansini, Fauzia Albertin

Cineca

Il progetto MIC Digital Library



Infrastruttura

I.PaC, l'*Infrastruttura e servizi digitali per il Patrimonio Culturale*, è lo **spazio dei dati** progettato per conservare, gestire e arricchire il patrimonio culturale digitale del Paese, in linea con le principali strategie nazionali ed europee. Nasce dall'esigenza di superare la frammentarietà dei sistemi di fruizione e dal bisogno di gestire dati stratificati ed eterogenei per formato, tipologia, dominio di appartenenza e politiche di protezione, secondo modelli concettuali flessibili e in sicurezza.

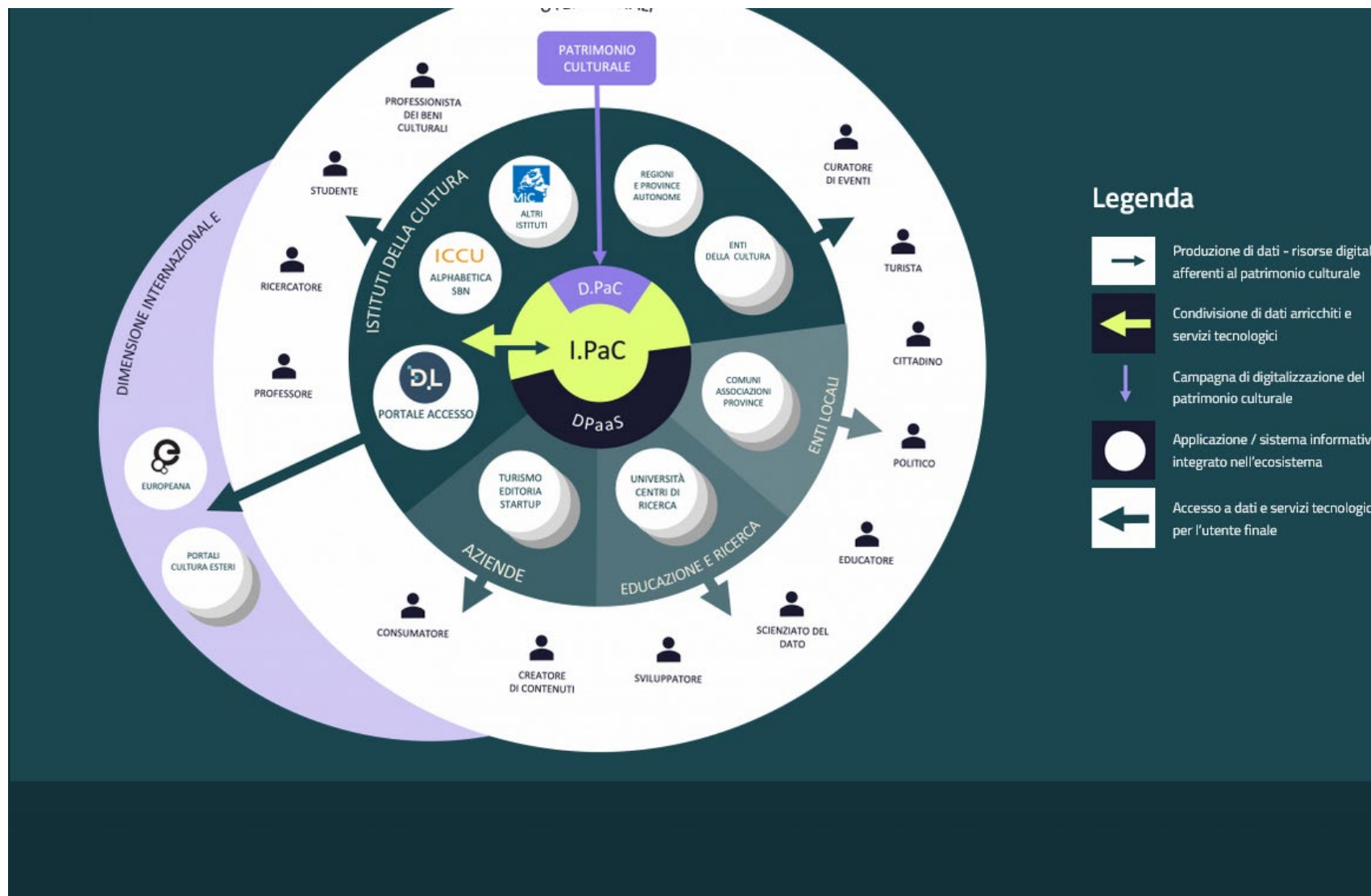
I.PaC presenta un **complesso sistema di servizi digitali avanzati**, basati su tecnologie innovative orientate al *Cloud*. Implementa funzioni relative alla **gestione e all'arricchimento delle risorse digitali**, fondate sia su modelli e schemi predefiniti (motori a regole e ontologie) sia su algoritmi di intelligenza artificiale (AI), ed **espone un ampio catalogo di API di cooperazione applicativa** (in lettura e scrittura) relative a dati di dominio e cross-dominio.

Scenario di riferimento

Lo sviluppo di I.PaC, e più in generale di un ecosistema digitale del patrimonio culturale italiano, rientra tra le azioni strategiche delineate nel *Piano nazionale di digitalizzazione del patrimonio culturale*, e si inserisce all'interno del più ampio progetto di trasformazione digitale promosso dalla Digital Library del Ministero della Cultura per il quinquennio 2022-2026.



Il progetto MIC Digital Library



La macchina delle cartucce di AI

The screenshot displays the 'AI SERVICES DEMO' web interface. At the top, there is a 'Login' button and a settings icon. The main heading is 'AI Services Demo', accompanied by a stylized 'AI' logo with circuit connections. Below this, the interface is divided into two sections: 'Available services' and 'Unreleased services'. The 'Available services' section contains 14 green cards, each representing a different AI service with its name and a brief description. The 'Unreleased services' section contains one yellow card for 'Italia9B'. At the bottom, it states 'Provided by Cineca'.

Available services						
Gliner Zero shot NER	Keras Keras OCR for text detection and recognition.	Kosmos2 Image caption generator	LayoutParser Layout recognition for magazines and newspaper	Llama31 Llama 3.1 8B LLM	NERs NERs - named entity recognition	Oemer Optical music recognition
P2PaLA Handwritten lines recognition	Qwen2 Multimodal LLM	RealESRGan General purpose image upscaler	StableDiffusion Text to image model	Tagging Text tagging	Tesseract Multilingual OCR for printed text	TrOCR Handwritten text recognition (single line)
Whisper Italian speech to text		Yolo Object detection and recognition				

Unreleased services
Italia9B First iGenius foundational Italian LLM

Provided by Cineca

<https://aiws.hpc.cineca.it/>

Layout Parser

Layout Parser Workflow v 0.5



active projects:

RDC_1939_09

RDC_1939_10

RDC_1939_11

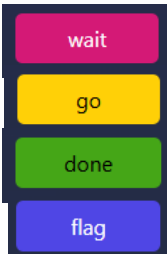
<https://aiws.hpc.cineca.it/>

Istruzioni

Scopo della piattaforma è allenare il motore AI di Layout Parser ad identificare correttamente le differenti categorie di testo presenti all'interno di una pagina di quotidiano – **titolo, testo, fotografia, disegno e pubblicità**

Una volta entrati nella cartella di progetto ci troveremo davanti le varie pagine di quotidiano identificate da un codice e di colore differente

- Il codice identifica in modo univoco la pagine
- Il colore identifica le azioni che sono già state fatte su di essa:



La pagina può esser messa in lavorazione

Il lavoro è stato completato

La pagina da attenzionare in quanto vi sono problematiche o dubbi

A screenshot of a web application interface. At the top, it says 'LPW / P3' and 'logged as: Fauzia'. The main title is 'RDC_1939_11'. Below it, the section 'Pages' contains a grid of 150 small colored boxes, each labeled with a page ID from P3_001 to P3_150. The boxes are color-coded: yellow for 'go', green for 'done', and blue for 'flag'. Some boxes have small icons (a red 'X' or a blue 'Z'). At the bottom, a 'Colors legend' shows four colored boxes: pink for 'wait', yellow for 'go', green for 'done', and blue for 'flag'.

Istruzioni

Una volta entrati nella pagina che dobbiamo lavorare ci troveremo:

- l'immagine della pagina
- già evidenziate le aree identificate dal motore Layout Parser

Come si vede dall'immagine queste non sono corrette in quanto parte dei titoli e dei testi non è stato identificato in alcun modo

Scopo di questo lavoro è identificare **TUTTE** le aree della pagina secondo 5 categorie

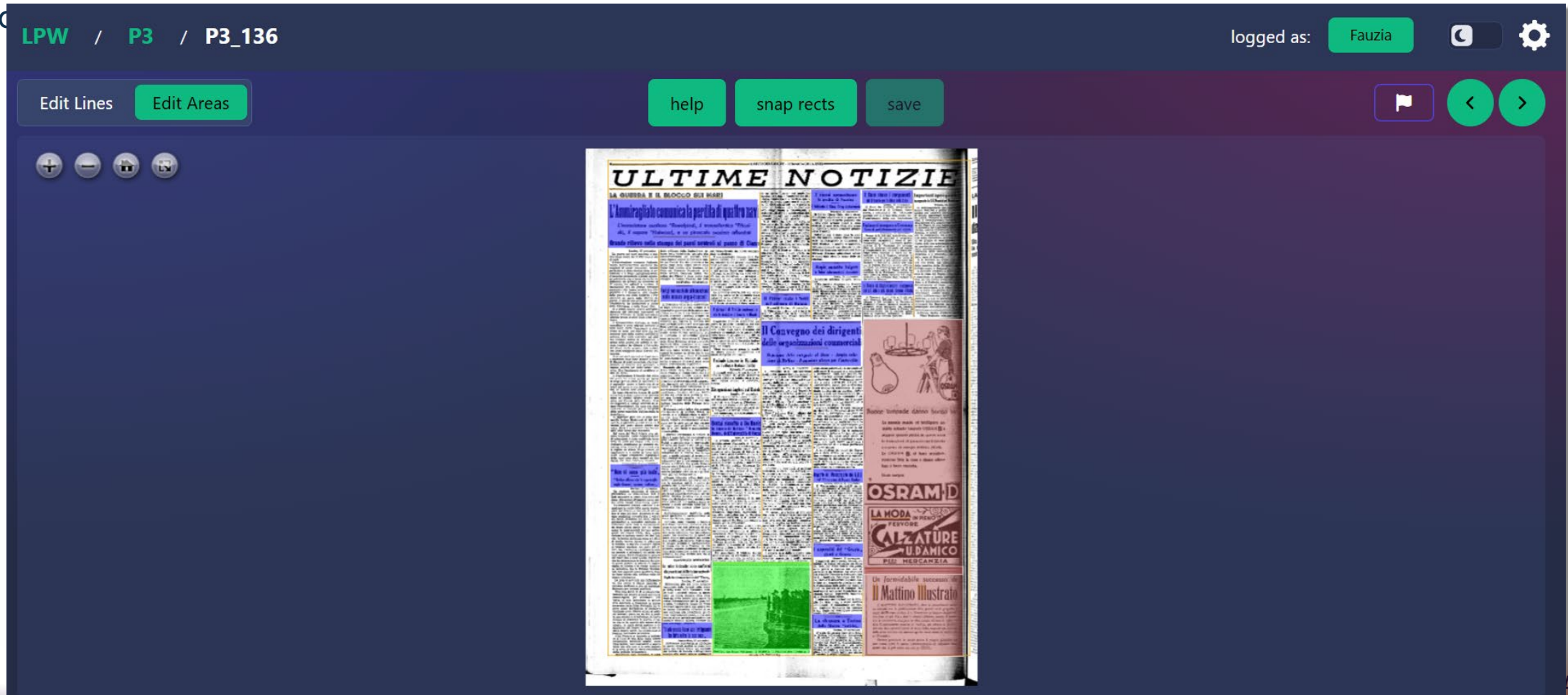
- 1 Titolo
- 2 Testo
- 3 Fotografia
- 4 Disegno
- 5 Pubblicità



Istruzioni

Nella parte superiore della pagina troviamo il nome della pagina che stiamo lavorando, il nostro user name e il pulsante delle impostazioni – in troveremo anche la corrispondenza colore-tipologia

La prima modalità di lavoro è **Edit Areas** – che ci permette di aggiungere, togliere o modificare le aree



Istruzioni

Posizionando il mouse sull'area di interesse e contemporaneamente uno dei numeri da 1 a 5 il sistema identificherà automaticamente l'area assegnandogli la categoria

1	Titolo
2	Testo
3	Fotografia
4	Disegno
5	Pubblicità

La dimensione dell'area identificata automaticamente viene determinata dalle linee guida – evidenziate in giallo.

Nel caso in cui l'area identificata automaticamente non fosse corretta, questa va modificata:

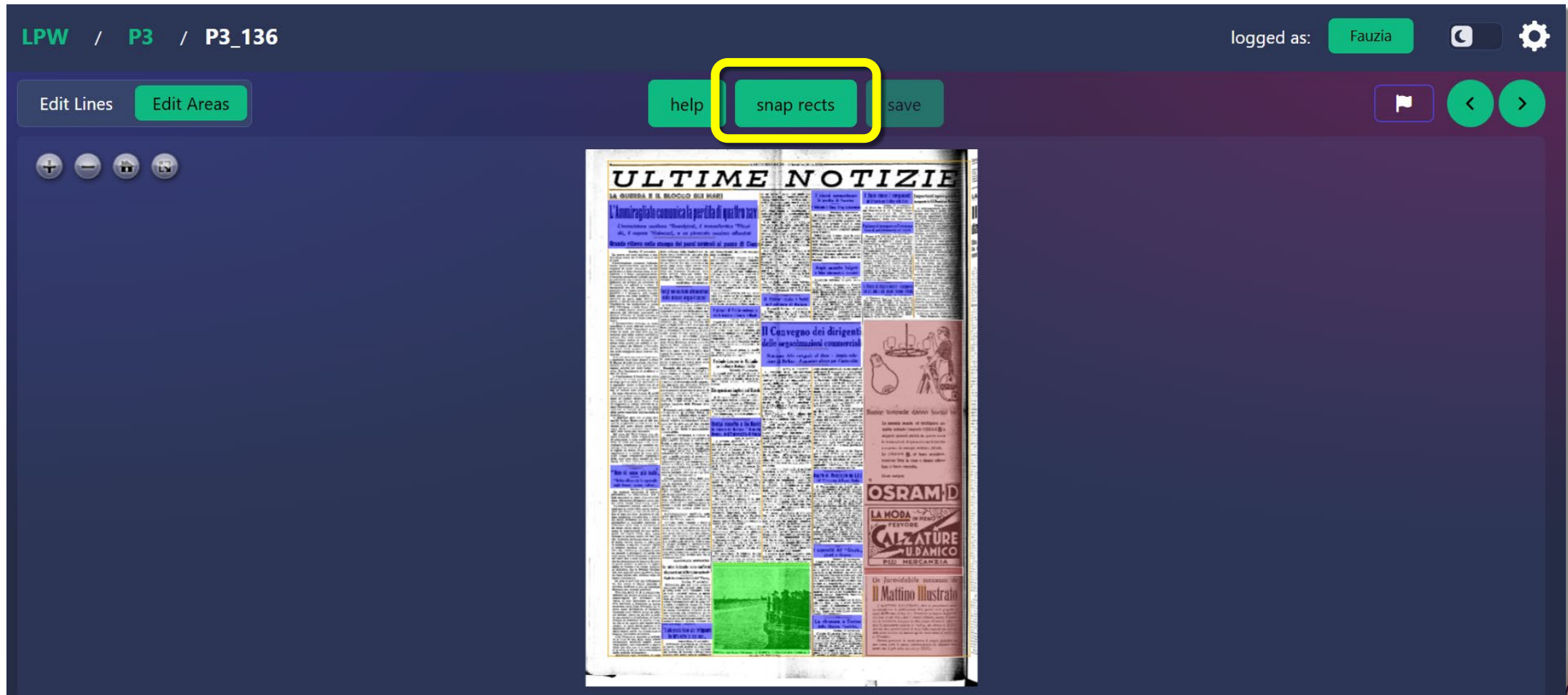
- agendo sui cursori ed allargandola o stringendola sino alla linea guida successiva.
- tenendo contemporaneamente premuto SHIFT è possibile muoversi al di fuori delle linee guida

E' possibile agire in ogni momento sulle aree identificate e compiere ulteriori azioni successive quali:

- **Split:** posso separare un area in due premendo S per lo split verticale o D per quello orizzontale. Apparirà una linea di divisione mobile che verrà confermata rilasciando il mouse
- **Merge:** selezionando due aree e premendo J le aree vengono unite in una unica

Istruzioni

Nel caso in cui si riscontrassero errori nelle linee guida PRIMA di cominciare il lavoro sulle aree è possibile chiedere al sistema di rettificare le linee agendo tramite il pulsante Snap rects



Istruzioni

Se gli errori sussistono è possibile entrare nella seconda modalità di lavoro: **Edit Lines**

In questa modalità è possibile:

- Spostare i marker di inizio e fine linea cliccando e trascinandoli
- Cancellarli cliccando e premendo CANC
- Creare nuovi marker cliccando e premendo CTRL

The screenshot shows the LPW software interface in 'Edit Lines' mode. The top navigation bar displays 'LPW / P3 / P3_136' on the left, 'logged as: Fauzia' in the center, and a settings gear icon on the right. Below the navigation bar is a toolbar containing 'Edit Lines' and 'Edit Areas' buttons, 'help' and 'save' buttons, and navigation arrows. The main workspace displays a newspaper page titled 'ULTIME NOTIZIE' with yellow markers highlighting specific lines of text. A small toolbar with zoom and pan icons is visible on the left side of the workspace.

Alcune norme

Titolo



- I titoli vanno suddivisi se provvisti di titolo e sottotitolo



- in questo caso no, si è semplicemente andati a capo



- anche la titolazione dei paragrafi è titolo

Alcune norme

Titolo



- Anche le pubblicità hanno i loro titoli



- in caso di decorazioni tipografiche (come linee o decorazioni) cerchiamo di non includerle

IL RESTO DEL CARLINO - 2 Novembre 1939 A. XVIII

Alcune norme

Testo

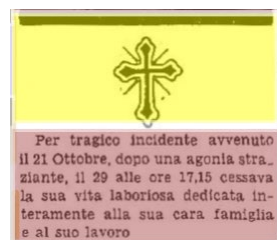
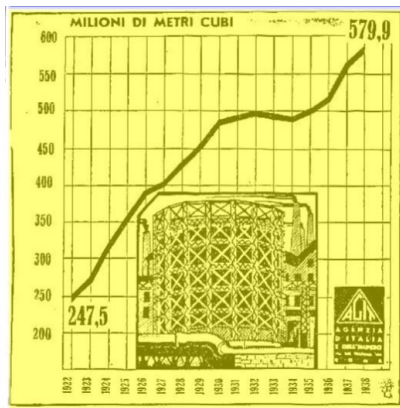
BOLLETTINO METEOROLOGICO					
CITTA	Tend. barom.	Stato del cielo	Stato del mare	Temperature	
				Massima	Minima
Bologna	variabile	schiarita	—	+ 8.7	+ 4.7
Roma	variabile	coperto	—	+ 15.0	+ 11.2
Atene	variabile	variabile	—	+ 18.8	+ 5.7
Torino	variabile	variabile	—	+ 11.2	+ 4.5
Genova	variabile	coperto	calmo	+ 14.6	+ 12.6
S. Remo	variabile	sereno	calmo	+ 17.4	+ 11.4
Venezia	variabile	schiarita	calmo	+ 8.0	+ 6.0
Trieste	variabile	schiarita	calmo	+ 12.8	+ 9.8
Trento	variabile	sereno	—	+ 8.3	+ 2.0
Modena	variabile	sereno	—	+ 10.0	+ 6.0
Firenze	variabile	coperto	—	+ 12.8	+ 7.8
Milano	variabile	sereno	calmo	+ 13.2	+ 4.8
Ancona	variabile	coperto	calmo	+ 15.0	+ 9.8
Napoli	variabile	coperto	calmo	+ 13.0	+ 11.0
Foggia	variabile	sereno	—	+ 16.5	+ 6.0
Bari	variabile	sereno	calmo	+ 16.5	+ 8.9
Lecco	variabile	sereno	—	+ 17.4	+ 8.0
Parma	variabile	sereno	calmo	+ 17.8	+ 8.4
Brescia	variabile	sereno	calmo	+ 17.0	+ 13.8
Palermo	—	—	—	—	—
Catania	variabile	sereno	calmo	+ 18.6	+ 5.0
Cagliari	variabile	sereno	calmo	+ 17.0	+ 7.4
Assisi	variabile	coperto	—	+ 17.5	+ 10.0
Frosinone	variabile	sereno	calmo	+ 22.2	+ 8.5
Reggio Emilia	variabile	coperto	calmo	+ 20.2	+ 17.2
Verona	variabile	sereno	calmo	+ 14.1	+ 12.2
Como	variabile	sereno	calmo	+ 18.8	+ 10.0

VENEZIA, 28 - PREVISIONI DEL TEMPO PER IL VENETO, L'EMILIA, LA ROMAGNA E LE MARCHE VALEVOLI FINO ALLE 18 DEL 29: Condizioni mediocri. Venti deboli occidentali. Cielo coperto. Nebbie e foschie diffuse. Mare poco mosso.

- le tabelle sono testo
- è buona norma selezionarle singolarmente separandole dal resto del testo

Alcune norme

Disegno



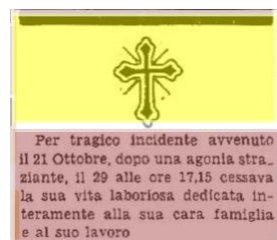
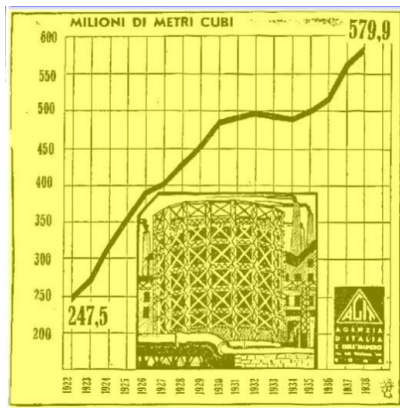
sera scende, dal piroscafi si spande sul mare il canto vibrante degli inni della Patria e della Rivoluzione.

Balbo inaugura in Cirenaica nuove importanti opere

- Grafici e mappe sono considerati disegno
- se dovesse esser presente una didascalia va inclusa nell'area
- anche marker tipografici come questo sono disegni
- se di questo tipo invece non vanno inclusi

Alcune norme

Disegno



sera scende, dal piroscifi si spande sul mare il canto vibrante degli inni della Patria e della Rivoluzione.

Balbo inaugura in Cirenaica nuove importanti opere

- Grafici e mappe sono considerati disegno
- se dovesse esser presente una didascalia va inclusa nell'area
- anche marker tipografici come questo sono disegni
- se di questo tipo invece non vanno inclusi

Alcune norme

Fotografia



- Vanno incluse anche le didascalie
- nel caso ci sia un'immagine con un titolo, il titolo va inserito come titolo, e non incorporato nell'immagine

Alcune norme

Publicità



- Le singole pubblicità vanno separate



- Disegni e testi all'interno delle pubblicità non vanno ulteriormente separati

Alcune norme

Publicità

Spesso la categorizzazione pubblicità o testo non è immediata. La logica che si è scelto di seguire è la seguente: se è l'annuncio è stato pagato per esser pubblicato è da considerarsi pubblicità

Lavori o immobili

AVVISI D'INDOLE COMMERCIALE
L. 2.50 per parola

A pagamento mensile: Impermeabili, paltò, novità, stoffe. Gianni, D'Azeglio 46. 12524

ACQUISTIAMO autocarri, automobili, macchinario ogni specie fuori uso. Prezzi massimi. Pasquali, Oriani 40. Telefono 20983. Bologna. 12620

ADDEZIONATRICI, calcolatrici, macchine da scrivere. Vastissimo assortimento. Cambi. Fornisconsi rivenditori. U.M.A., Telefono 33-666, Farini 14 interno. 12315

APRILIA acquistata giugno carrozzeria speciale come nuova vendesi L. 31500. Scrivere CASSETTA 7 V UNIONE PUBBLICITA' ITALIANA. Bologna.

BIANCHI S 9 Cabriolet lussuoso perfetto adattissimo metano svendo 13.500. Acquisto Ford 3 litri oppure 522 purchè perfettissima. Cassetta 33822 Z Unione Pubblicità Italiana. Trieste.

CAMION Lancia 35 quintali 9750. Berlino, Alfa, Ballila, 514. 509, vere occasioni. Martinelli, S. Giorgio 3. 12636

CAMIONISTI acquistate gassogeno «Dux» vera economia 0,80 carbonella sostituisce un litro carburante qualsiasi applicazione. Martinelli, S. Giorgio 3. 12638

FIAT 500 1100 1500 Lancia Aprilia non oltre otto mesi vita acquistansi. Indicare prezzo ristretto contanti. Scrivere CASSETTA 6 T UNIONE PUBBLICITA' ITALIANA. Bologna. 12564

Nuove nascite

FIOCCHI BIANCHI

△ GIAN CARLO VACCARI annuncia con gioia la nascita del fratellino

GIAN FRANCO

Bologna, 22 novembre 1939-XVIII.
Via Savioli 10.

Necrologi

Dopo lunga e penosa malattia, munita dei conforti religiosi, si è spenta a 55 anni

Chiara Cavallini
in Giordani

Addoloratissimi ne danno il triste annuncio il marito Dottor GIAMBATTISTA, i figli LUIGI e MARIA CONCETTA, la madre TERESA PIATEGI ved. CAVALLINI, i fratelli PIETRO e LUISA col marito Rag. SALVATORE FALCHI, la cognata MARIA GIORDANI, i nipoti e parenti tutti.

Non fiori, ma opere benefiche e preghiere.

Bologna, Via Galliera 12.
Lugo, Via Garibaldi 38.

Mercoledì 29 corr. alle ore 9 nella Basilica di S. Maria Maggiore in Via Galliera (presente la salma) sarà celebrata una Messa; pocca la cara Salma sarà trasportata a Lugo dove alle ore 11, dopo rinnovate esequie nel Santuario del Molino, verrà tumulata nel sepolcro di famiglia.

Film e spettacoli

TRAPPOLA D'AMORE

Oggi al MODERNISSIMO

E' un film spassoso diretto da Tarazzo con Carla Candiani, Lilla Landini, Osvaldo Valenti, Giuseppe Pirelli - Prima visione.

In questo caso i cinema non hanno pagato per pubblicare, è quindi da considerarsi testo

Spettacoli d'oggi

MANZONI - «I forzati della libertà» Tylor, Sc. C. Marimba Mastrolia

MOD. - «Trappola d'amore» G. Pirelli, Osvaldo Valenti, C. Candiani, Lilla Landini

FULGOR - «Eravamo 7 vedovi» Anna Gandusio, Laura Nucci, M. D'Amico

CENTRALE - «Fornaretto di Venezia» Don. «Fratello Ballochino» Fr. Maria

IMPERIALE - «Assenza inaspettata» Alida Valli, Amedeo Nazzari, La Voce

VERDI - «Eroe sconosciuto» L. Landini «Voglio essere anata» Douglas

CONTRAVALLI - «Carnet di ballo» Richard «Tesoro del Parione» E. Pirelli

APOLLO - Comp. L'aereo della comita. Sch. «Napoli che non muore» G. Pirelli

MARCONI - «Dolce inganno» Enrico Tone «Miasma pericoloso» Fr. Maria

CARLUCCI - «Katia» enorme successo ultimo giorno e «Vorrei volare»

DUSE - Var. Comp. Riv. Carovana Anna. sch. «Belle e brutte» si sposterà

REX - «Aspetto una signora» Lilla Landini «Cardinale Richelieu» G. Pirelli

ROMA - «Bel oca di permesso» T. Pirelli

OLIMPIA - «Una notte d'oliva» L. Landini

MODERNO «La Signora dalle Ombre»

SAVOIA - «Le avventure di Tom Sawyer» con Tommy Kelly, T. Pirelli

MEDICA - Ore 15. Letture e spettacoli di Eugénia e rivista «Follie»

Alcuni casi peculiari

PICCOLI AVVISI
MINIMO 10 PAROLE OGNI AVVISO

Si ricevono presso la
UNIONE PUBBLICITA' ITALIANA

N.B. - Tutti gli avvisi provenienti da
agenzie sono soggetti alla tariffa « Com-
merciale ».

Alcuni casi peculiari

Anno LV N. 284 - Italia Impero Colonia, cent. 30

il Resto del Carlino

Bologna - Giovedì 30 Novembre 1939-XVIII

ABBONAMENTI
ITALIA IMPERO COLONIE, Anno L. 75 Sem. L. 38 Trim. L. 20
Cin. Fed. del Nord Anno L. 87 Semestre L. 44 Trimestre L. 23
D.E.F. L'ESTERO, Anno L. 140 Semestre L. 41 Trimestre L. 45
Numero annuo L. 630 - Direzione e Amm. SOCIOFA, Via S. Luigi N. 5 -
Tel. 2111 - 2112 - 2113 - 2114 - 2115 - 2116 - 2117 - 2118 - 2119 - 2120 - 2121 -
Inserzioni via Grafica - Spedite in abbonamento postale
C. C. postale n. 6-747

TARIFFA PER LE INSERZIONI
Freti per ann. di abbas. (richiesta di una colonna) Piacenza L. 9 -
Commerciale L. 6 - Mortara L. 5 - Cronaca L. 10 (includo
20mm.) Piacenza Avv. vedi tariffe in testa alle varie rubriche
Pagamento anticipato - Tassa sulla pubblicità 10-14 p. l. n. 26-903
esclusivamente a SOCIOFA, Via S. Luigi n. 5 - tel. 26-903
UNIONE PUBBLICITA' ITALIANA S. A.

il Resto del Carlino

ABBONAMENTI
per il 1940-XVIII-XIX

ITALIA IMPERO COLONIE	Ann.	Sem.	Trim.
Del lunedì settiman.	75-	38-	20-
Con Fed. de lunedì	87-	44-	23-

ESTERO

Del lunedì settiman.	160-	81-	41-
Con Fed. de lunedì	186-	94-	48-

EDIZIONE DELLA SERA
il Resto del Carlino
Anno L. 75 Sem. L. 38 Trim. L. 20

ABBONAMENTI CUMULATIVI
IL RESTO DEL CARLINO «1»
La Rivista del «Popolo d'Italia»
Pubblicazione mensile dei principali avvenimenti della Patria, dell'Impero, dell'Avia e dello Sport - L. 180
Rivista «Rassegne mensili della Repubblica» - L. 100
Rivista «Fascista» 1940-XVIII-XIX - Valore trimestrale (Rivista edita da il Popolo d'Italia) - L. 80
L'Illustrazione Italiana - L. 200
Scienze - Quotidiano di divulgazione di Scienze teoriche e Arte applicata - L. 100
Dress - Rivista mensile di Moda e modista - L. 100
Rivista - Antologia letteraria di Letteratura italiana - L. 100
Giornale - Grande rivista mensile di Lettere, Scienze, Arte e di Piacenza società fascista - L. 100
La Donna Italiana - Rivista mensile di Lettere, Scienze, Arte e di Piacenza società fascista - L. 100
Civiltà - Rivista mensile - L. 100
I Paesi del Mondo - Sotto la Direzione Società Geografica Italiana. Rivista mensile del 1910 gruppo di tutti il mondo - L. 80
Calendario Fascista De Agostini 1940 - Circa 500 copie di

Nel vostro esclusivo interesse citate sempre nei vostri ordini e offerte
"IL RESTO DEL CARLINO"

- Queste sono pubblicità in quanto riportano costi e tariffe per gli abbonamenti al giornale

Alcuni casi peculiari

STATO CIVILE DI BOLOGNA

Denuncia del 25 Novembre 1939-XVIII

Nati	6
Morti	16
Matrimoni	..

FIOCHI BIANCHI

△ GIAN CARLO VACCARI annun-
cia con gioia la nascita del fratellino

GIAN FRANCO

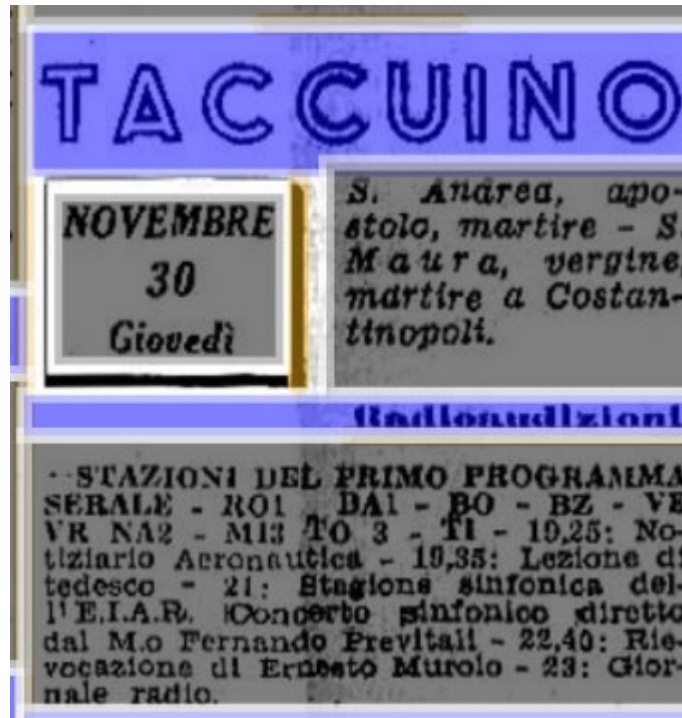
Bologna, 22 novembre 1939-XVIII.
Via Savioi 10.

Al nati annunciati in
questa rubrica l'Unio-
ne Pubblicità Italiana
regala un biglietto del-
la Lotteria «E 42»
col quale possono vin-
cere molti milioni

Al piccolo Gian Franco
Vaccari è stato assegnato
il biglietto della Lotteria
E. 42 N. 40952 Serie I.

- E' un testo in quanto un dato sullo stato civile
- E' una pubblicità in quanto annuncio pagato dalla famiglia
- E' un testo in quanto riporta una notizia

Alcuni casi peculiari



- La data è testo – da non includere il riquadro

Immagini e titolo

Nel caso di un'immagine con didascalia ma anche titolo, che non ha corrispondenza con altre aree di testo, il titolo non va incorporato nell'immagine ma marcato come titolo



Titoli e sottotitoli non riunibili



In casi come questo si è deciso per ora di lasciare la suddivisione delle bounding box

Casi da attenzionare



Nel caso in cui la pagina sia stata acquisita leggermente storta non sarà possibile identificare correttamente le colonne, e il problema si porrà maggiormente nelle colonne a tutta altezza. Non vi è soluzione (al momento) per questa problematica.

Infatti:

- Layout Parser lavora solo su aree rettangolari, non è possibile creare trapezi per seguire la colonne non perfettamente verticali
- Suddividere la colonna in n rettangoli successivi andrebbe a confondere il modello

Soluzione: *fai come faresti** e lascia un flag in maniera tale che la pagina possa esser facilmente riconoscibile come «problematica»

* cerca di trovare il miglior compromesso tra aree incluse e aree sovrapposte

Questioni aperte

- Il sistema riconosce anche le gerarchie Titolo – Articolo?
 - Se si, lo stiamo allenando anche per questo?
 - Se no, dovremmo ri allenarlo?
 - O servono ulteriori motori (ancora da identificare)?
- [SIL] Layout Parser non crea le gerarchie
la nostra intenzione sarebbe di estrapolarla a valle
con del codice nostro ad hoc

Grazie dell'attenzione

<http://visitlab.cineca.it/>
visitlab@cinca.it