

Page Segmentation in Urdu Nastalique OCR

Qasim Ali

qasim.ali@kics.edu.pk

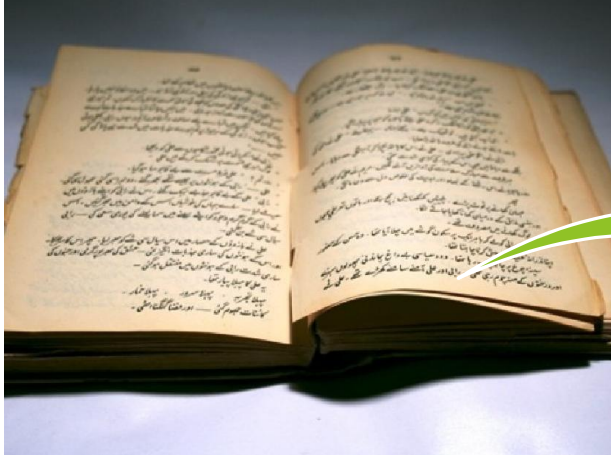
29-September, 2014

مرکز تحقیقات لسانیات



Center for Language Engineering
Al- Khwarizmi Institute of Computer Sciences
University of Engineering and Technology, Lahore, Pakistan





Optical Character Recognition

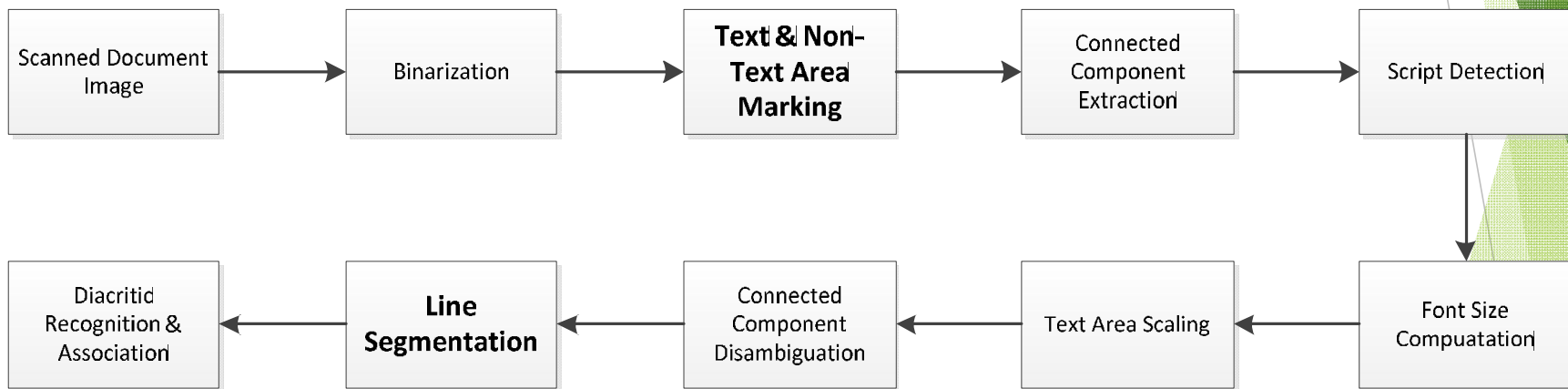


106

آئے ہوئے تمام لوگوں کا شکریہ ادا کرتا ہے اور مرحوم یا مرحومہ کی طرف سے معافی کا طلبگار ہوتا ہے کہ حیات میں ان سے کوئی غلطی ہوئی ہو، وہ معاف کر دیں اور خدا سے دعا کرو۔ پھر مجلس برخواست ہوتی ہے۔ سارے لوگ اپنے اپنے گھروں کی طرف چلے جاتے ہیں اہل ہمسائے اور برادری کے لوگ اس موت کے سر چلے جاتے ہیں 'جہاں پر ہمسائے ان کے لیے پائے اور اس کے ساتھ روٹیاں لے آتے ہیں 'وہ پیش کی جاتی ہیں۔ وہاں پھر فاتحہ پڑھی جاتی ہے۔ اس دوران یہ سارے رسوم خلیفہ کے ذریعے ادا ہوتے ہیں 'جہاں پر خلیفہ ہر وقت ان کے درمیان موجود رہتا ہے اور جو لوگ فاتحہ خوانی کے لیے آتے ہیں ان کے لیے فاتحہ خوانی کرتا ہے۔

Phases In OCR

1. Pre-processing



2. Classification and recognition

3. Post-processing

Nastalique Writing Style Complexities

Nastalique

پاک کلام

~~چمچ چمچہ جمہوریہ نستعلیق~~

~~چمچ / جمہوریہ / نستعلیق~~

نانک

تنبہ سنت

Bottom Up Approach

9

۱۰۰ ﴿اے رسول﴾! ہمیں کہہ دیجئے (حشم بصیرت لے کر) حل بھر کر اہم گزشتہ کے ادمار و
۱۰۱ زوال کے کھنڈرات دیکھو اور بحرموں کے انجام سے (درس عبرت لو)
۱۰۲ ﴿سورہ اہل 27: ایت 69﴾
۱۰۳ ماضی میں جھانکنے یا کھوج لگانے کے لیے علم آثار قدیمہ نے جو کہ اب ایک اہم سائنس
۱۰۴ ہے، نہایت اہم کام کیا ہے۔ علم آثار قدیمہ کو ایک جذباتی سائنس کہا جاتا ہے کیونکہ اس میں
۱۰۵ بے آرزو بیجان اور بہت زیادہ تکلیف دہ سائنسی طریقے شامل ہیں۔ پوری دنیا میں یہ سائنس
۱۰۶ اٹھارہویں صدی میں یوم پی آئی (Pompei) کی کھدائی کے ساتھ معرض وجود میں آئی۔
۱۰۷ پوم پی آئی شہر اٹلی کے جنوب مغرب میں واقع تھا۔ جو 79ء میں وسوویس
۱۰۸ (Visuvius) پہاڑ کے پھٹنے سے آتش فشاں کے طے کے سچے دہن سو گیا تھا۔ 90 فی صد شہر
۱۰۹ کی آبادی وقت میں باہر آگئی تھی لیکن 2000 لوگوں نے شہر چھوڑنے سے انکار کر دیا تھا۔ ہلکی
۱۱۰ مسام دار چٹان جسے پومک پتھر کہتے ہیں، اس کے چادلوں کے داؤوں کے برابر اور دوسرے
۱۱۱ مشتمل بھر جسامت کے 15 سنی میٹر (6 انچ) فی گھنٹہ کی رفتار سے بارش کی طرح گلیوں میں
۱۱۲ گرے اور پھر اسی پہاڑ سے ایک مادل خاک کس اور چٹان پر مشتمل آبا اور اس دم گھسے
۱۱۳ والے میٹیریل کے طوفان سے ہائی ماندہ پوم پی آئی کے لوگ یا باشندوں میں سے کوئی نہ بچا۔
۱۱۴ 19 گھنٹوں کے اندر یہ شہر 25 فٹ (8 مسر) موٹی طے کے طے کے سچے دہن ہو چکا تھا جو
۱۱۵ سینٹ کی طرح سخت ہو گئی۔ اس دہن شدہ شہر سے 1700 سال تک کسی نے کوئی تعلق نہ رکھا،
۱۱۶ پھر 1863 میں ماہرین آثار قدیمہ نے پوم پی آئی کی کھدائی شروع کی۔
۱۱۷ علم آثار قدیمہ کی ابتدائی ٹیوٹوٹا یا فروغ عہد جدید کی طلوع سحر کے ساتھ فنون لطیفہ
۱۱۸ کے ماہرین کی رومانوی تحریک مستطین ہوئی جو پہلے ہی جذبے اور تصور کی گہری سوچ میں غرق
۱۱۹ تھے، ان کو شجر و بیابان تصویر کی مانند اور خوشنما مقامات کی طرف لایا گیا جن میں کھنڈرات بھی
۱۲۰ شامل تھے اس سے پہلے تو کوئی بھی اس کا استعمال نہ کرتا تھا بلکہ تاریخ کے اس طے میں کوئی
۱۲۱ دلچسپی نہ لیتا تھا۔ علم آثار قدیمہ کو جس کا تعلق قدیم زمانے کے مطالعے سے ہے ابتدائی زندگی
۱۲۲ میں — کس سے جو بھی شہادت پائی جاسکتی تھی حاصل کی۔ درحقیقت سائنس اور رومانس کی
۱۲۳ دوہری مہک حاصل کی، جو نئی لوگوں نے اپنے ماضی کے بارے میں کھوج لگانے شروع کی اور
۱۲۴ مستقبل کے بارے میں سراخ لگانے کی امید میں اپنے ماضی کے بارے میں سوال کیا، لو اس

Challenges In Two Column Segmentation

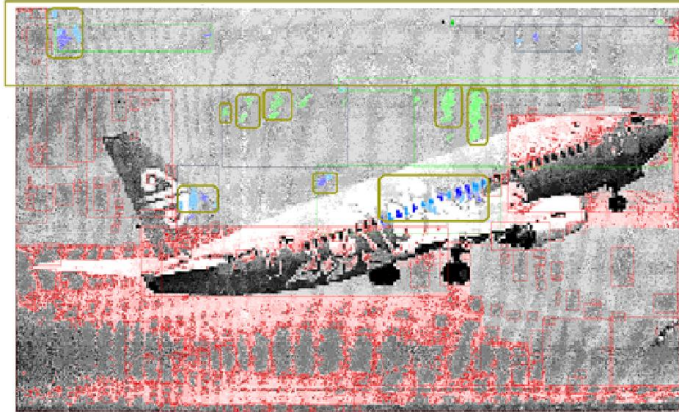
”لکھے ہیں“ آپا ہے۔

پے	رئج	ہدم	دریں	سال	سی
عم	قلدہ	کردم	ہ	ایں	ہاری
سو	اردو	کی	آراسلا	کر	قہاں

لکھے اس نے ہندی زبان کو دیا
دیا نظم اردو کو یہ مرثیہ

جائیں۔ بہر حال یہ ایک حقیقت ہے کہ گائٹم اور نوو ویلیوف کی دریافت

بالا بینڈ سٹرکچر کی تصدیق کر لی تو اس سے مزید تجربات کی راہ کھلی، جن میں سے ایک یہ تھا کہ گرافین کمرے کے درجہ حرارت پر بھی کوآٹم ہال اثر

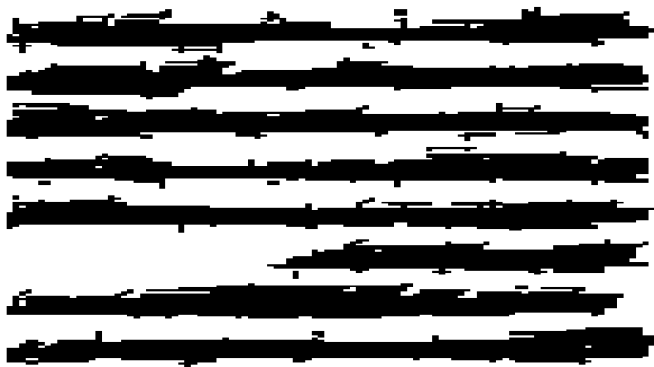


Layout Extraction

دعا کریں۔ پھر مجلس برخواست ہوتی ہے۔ سارے لوگ اپنے اپنے گھروں کی طرف چلے جاتے ہیں البتہ ہمسائے اور برادری کے لوگ اس موت کے گھر پر چلے جاتے ہیں، جہاں پر ہمسائے ان کے لیے چائے اور اس کے ساتھ روٹیاں لے آتے ہیں، وہ پیش کی جاتی ہیں۔ وہاں پھر فاتحہ خوانی ہوتی ہے۔ اس دوران یہ سارے رسوم خلیفہ کے ذریعے ادا ہوتے ہیں، جہاں پر خلیفہ ہر وقت ان کے درمیان موجود رہتا ہے اور جو لوگ فاتحہ خوانی کے لیے آتے ہیں ملن کے لیے فاتحہ خوانی کراتا ہے۔

اس کے بعد تین دن تک اس گھر میں کچھ بھی نہیں پکتا بلکہ تین دن تک ہمسائے کے لوگ، قبیلے کی خواتین اور رشتہ دار روزانہ تین اوقات میں چائے اور سالن دروٹیاں

Scanned image of font size 14

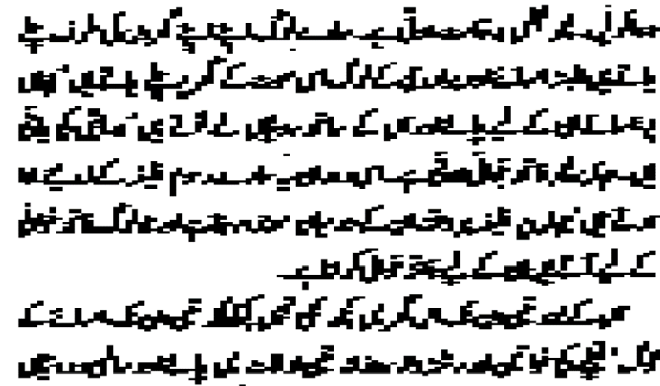


Resultant image after smearing

دعا کریں۔ پھر مجلس برخواست ہوتی ہے۔ سارے لوگ اپنے اپنے گھروں کی طرف چلے جاتے ہیں البتہ ہمسائے اور برادری کے لوگ اس موت کے گھر پر چلے جاتے ہیں، جہاں پر ہمسائے ان کے لیے چائے اور اس کے ساتھ روٹیاں لے آتے ہیں، وہ پیش کی جاتی ہیں۔ وہاں پھر فاتحہ خوانی ہوتی ہے۔ اس دوران یہ سارے رسوم خلیفہ کے ذریعے ادا ہوتے ہیں، جہاں پر خلیفہ ہر وقت ان کے درمیان موجود رہتا ہے اور جو لوگ فاتحہ خوانی کے لیے آتے ہیں ملن کے لیے فاتحہ خوانی کراتا ہے۔

اس کے بعد تین دن تک اس گھر میں کچھ بھی نہیں پکتا بلکہ تین دن تک ہمسائے کے لوگ، قبیلے کی خواتین اور رشتہ دار روزانہ تین اوقات میں چائے اور سالن دروٹیاں

Binarized version

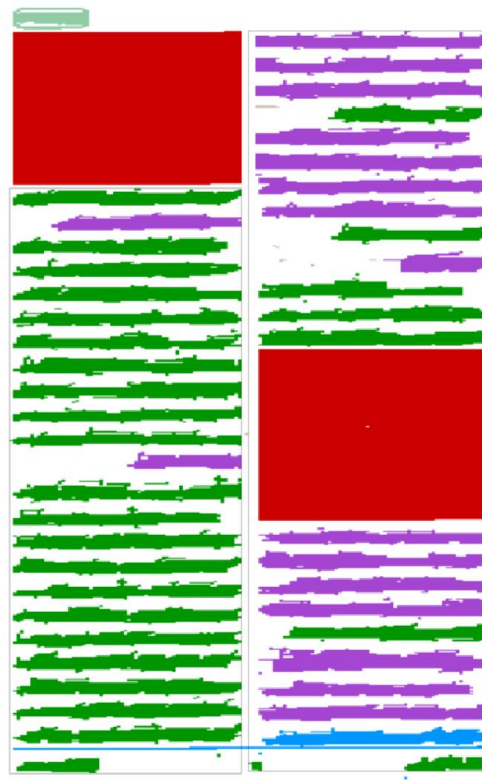


725x1306 pixels down sampled to 182x109

Layout Extraction



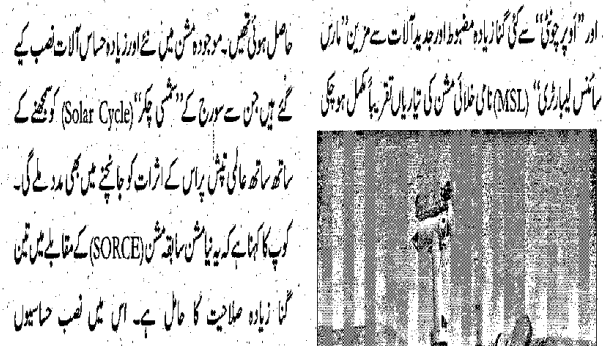
Page with single column



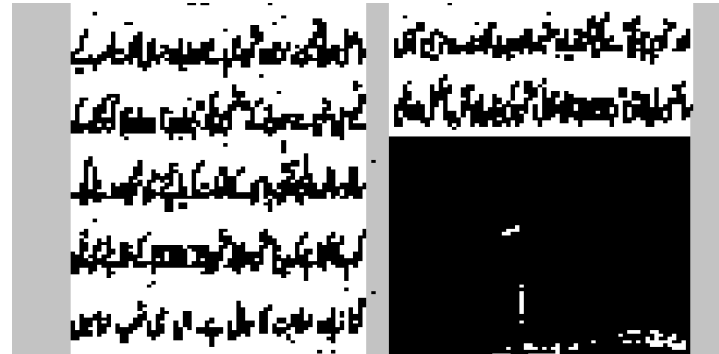
Page with two columns



Layout Extraction



Page with two columns



White obstacle shown in grey color

Layout Extraction

آئے ہوئے تمام لوگوں کا شکر یہ ادا کرتا ہے اور مرحوم یا مرحومہ کی طرف سے معافی کا طلبگار ہوتا ہے کہ اگر حیات میں ان سے کوئی غلطی ہوئی ہو، وہ معاف کر دیں اور خدا سے دعا کریں۔ پھر مجلس برخواست ہوتی ہے۔ سارے لوگ اپنے اپنے گھروں کی طرف چلے جاتے ہیں البتہ ہمسائے اور برادری کے لوگ اس موت کے گھر پر چلے جاتے ہیں، جہاں پر ہمسائے ان کے لیے چائے اور اس کے ساتھ روٹیاں لے آتے ہیں، وہ پیش کی جاتی ہیں۔ وہاں پھر فاتحہ خوانی ہوتی ہے۔ اس دوران یہ سارے رسوم خلیفہ کے ذریعے ادا ہوتے ہیں، جہاں پر خلیفہ ہر وقت ان کے درمیان موجود رہتا ہے اور جو لوگ فاتحہ خوانی

Color coded output of a text area

106

آئے ہوئے تمام لوگوں کا شکر یہ ادا کرتا ہے اور مرحوم یا مرحومہ کی طرف سے معافی کا طلبگار ہوتا ہے کہ اگر حیات میں ان سے کوئی غلطی ہوئی ہو، وہ معاف کر دیں اور خدا سے دعا کریں۔ پھر مجلس برخواست ہوتی ہے۔ سارے لوگ اپنے اپنے گھروں کی طرف چلے جاتے ہیں البتہ ہمسائے اور برادری کے لوگ اس موت کے گھر پر چلے جاتے ہیں، جہاں پر ہمسائے ان کے لیے چائے اور اس کے ساتھ روٹیاں لے آتے ہیں، وہ پیش کی جاتی ہیں۔ وہاں پھر فاتحہ خوانی ہوتی ہے۔ اس دوران یہ سارے رسوم خلیفہ کے ذریعے ادا ہوتے ہیں، جہاں پر خلیفہ ہر وقت ان کے درمیان موجود رہتا ہے اور جو لوگ فاتحہ خوانی کے لیے آئے ہیں ان کے لیے فاتحہ خوانی کرتا ہے۔

اس کے بعد تین دن تک اس گھر میں کچھ بھی نہیں پکا بلکہ تین دن تک ہمسائے کے لوگ، قلیبی کی خواہن اور رشتہ دار روزانہ تین اوقات میں چائے اور سالن دروٹیاں لاتے ہیں۔ اس کے ساتھ ہی مرد لوگ بھی آس پاس سے اور دور دراز علاقوں سے تعزیت کے واسطے آتے رہتے ہیں۔ پھر ان کے سامنے چائے و روٹی سے خاطر تواضع بھی کی جاتی ہے۔ تین دن تک برادری کے لوگ بھی ساتھ رہتے ہیں اور سوگواروں کی ہر طرح سے دلجوئی کرتے رہتے ہیں۔ تیسرے روز کی شام کو ایک بھیڑ و ذبح کر کے نذر دینا زکات بند و بست کیا جاتا ہے جس میں برادری کے لوگوں کے علاوہ ہمسائے اور دیگر رشتہ دار شامل ہوتے ہیں۔ اس رسم سوگم کی ادا ہو گئی کے بعد بھی بعض خواہن اور مرد جو اس دوران نہیں آسکے ہیں، ایک ہفتے تک تعزیت کے لیے آتے رہتے ہیں اور مرحومین کے لیے فاتحہ خوانی کرتے ہیں۔

یہ بات بھی واضح رہے کہ اس رسم کی ادا ہو گئی میں امیر و غریب کا کوئی امتیاز نہیں۔ سب کے لیے ایک قسم کی رسم ہے اور سب کی دلجوئی کی جاتی ہے۔ اس طرح غمزدہ خاندان کو حوصلہ ملتا ہے اور ایک قسم کی بھائی چارگی کا احساس بڑھتا رہتا ہے۔ یہ روایت اب تک برابر جاری ہے اور اس وقت ہنزہ والے خواہن پاکستان کے کسی بھی شہر میں رہتے ہوں، اسی رسم کی ادا ہو گئی میں برابر کا رہتے ہیں اور اس طرح غمزدہ خاندان کو کہیں بھی بے چارگی یا تنہائی کا احساس نہیں ہوتا۔

Page with single column

Accuracy - Text Extraction

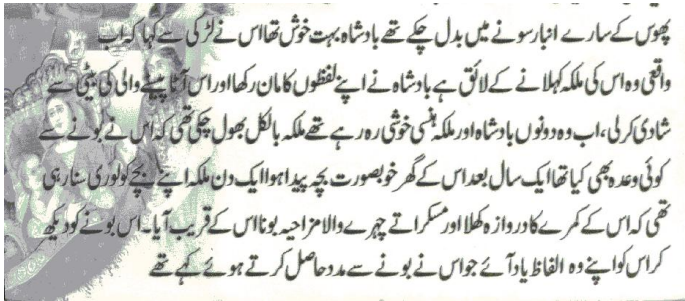
Total Images	1,257
Total Gold Textual CCs	2,092,782
Noise CCs	931,522
Gold Textual CCs excluding noise	1,161,260
Computed Textual CCs	1,140,959
Matched Textual CCs	1,140,712
Accuracy of Text as Text	98.23%

Accuracy - Line Segmentation

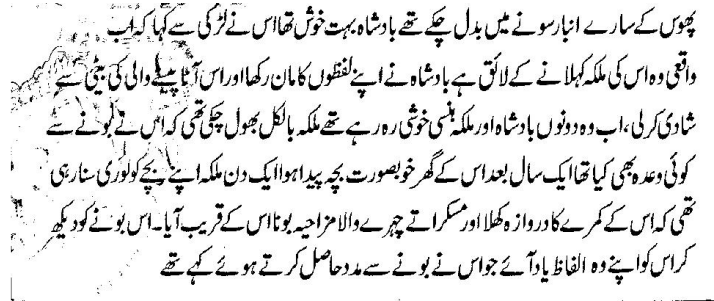
پر ہمسائے ان کے لیے چائے اور اس کے ساتھ روٹیاں لے آتے ہیں، وہ پیش کی جاتی

Font Size	Total Number of Lines	Correct Number of Lines	Accuracy
14	930	865	93.01%
16	820	758	92.44%
18	720	657	91.25%
20	570	439	77.02%
22	447	325	75.93%
24	700	659	94.14%
28	677	624	92.17%
32	775	679	87.61%
36	634	573	90.38%
40	519	443	85.36%

Layout Extraction Challenges



Half toner image of Font 22

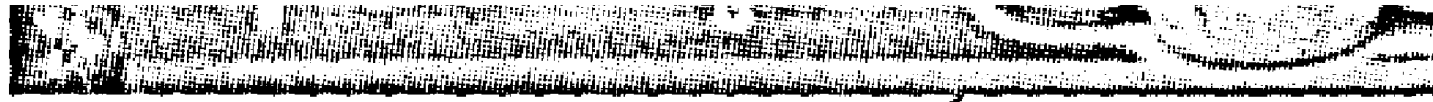


Binarized version



Smearred version

Layout Extraction Challenges



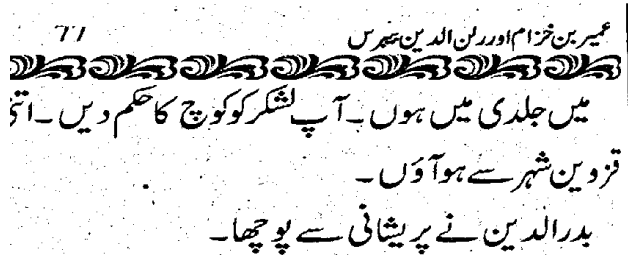
دو حصوں میں بٹ گئے۔ کچھ تو مکمل طور پر مسجد اور مدرسے تک محدود ہو کر رہ گئے، اور باقی ان سے

Main body joined with figure

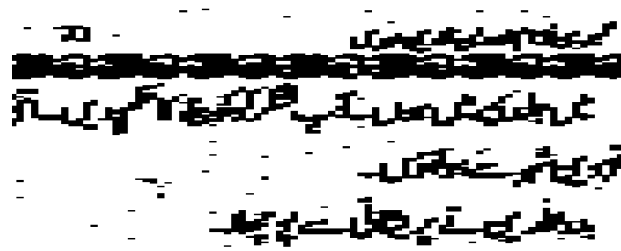


Resultant smeared image

Layout Extraction Challenges



Page Frame in Binarized version



No White Obstacle detected



Lines merged with page frame and detected as figure