



Распознавание таблиц в форматированных документах

Илья Копаничук, Инга Очнева, Александр Огальцов,
Мариам Каприелова, Евгений Финогеев, Александр Кильдяков,
Юрий Чехович



В рамках исследования требуется выделять таблицы из форматированных документов, чтобы или не учитывать их при поиске заимствований, или в дальнейшем строить индекс таблиц и искать в нем заимствования с учетом особенностей форматирования

- Решение предполагается мультиязычным
- Решение должно удовлетворять промышленным нагрузкам



Задана выборка:

$$\mathcal{D} = \{d_i, \mathbf{y}_i\}_{i=1}^m, \quad d_i = (x_1^i, \mathbf{m}_1^i), \dots, (x_n^i, \mathbf{m}_n^i), \quad \mathbf{y}_i = y_1^i, \dots, y_n^i,$$

где d_i рассматриваемый документ с форматированием, задается последовательностью токенов x_j^i и соответствующим им форматированием \mathbf{m}_j^i на основе метаданных документа.

Требуется построить отображение:

$$f: \mathbb{D} \rightarrow \mathbb{Y},$$

где \mathbb{D} документы с заданным форматированием для слов, а \mathbb{Y} последовательности меток класса для слов.



Признаковое описание:

- Признаковое описание слова документа задается на основе форматирования документа: высоте, ширине и положении слова на странице
- Мультиязычность признакового описания: используется только геометрия токенов и ключевые слова
- Метки класса: $text = 0$, $table = 1$
- Подписи таблиц относятся к тексту, но имеют особое значение
- Под документом рассматривается только одна страница – все страницы многостраничного документа обрабатываются независимо

Последовательность работы предлагаемого метода:

1. Выделение линий документа (сегментация слов на странице)
2. Классификация линий документа
3. Перевзвешивание классов линий и присвоение меток словам



- Синтетические данные:
 - 10000 документов (Python → LaTeX)
 - Рандомизированы почти все параметры сборки
 - Разметка идет по размерам шрифтов текста, таблиц и подписей
 - 90% документов содержат и текст, и таблицы
- Реальные данные:
 - 2008 документов на двух языках (ru, en)
 - Размечены собственноручно
 - Только 13% документов содержат таблицы
- Обучающая выборка:
 - 2804 документа
 - 764 тысячи токенов
 - Реальные + синтетические данные 1:1
 - 7.2% токенов – таблицы
- Валидационная выборка:
 - 757 документов
 - 241 тысяча токенов
 - Реальные + синтетические данные 2:1
 - 6.4% токенов – таблицы
- Тестовая выборка:
 - 602 документа
 - 145 тысяч токенов
 - 100% реальные данные
 - 6.5% токенов – таблицы



Примеры синтетических сгенерированных страниц с таблицами

Ru, 1 колонка

Прогрессом тесно связано перемена истинную гушат кушанью укажию общию одарить миротом инертности багратионов. Заменяю бродни динимый акордные поварывавшие уронт реваа сложившия возросшие повашим несомну замойки. Спадающим сарлемшия уверявшия пылавые возгнотвориди, ста изгнривавшия цастрождою рачениа акоса передидиненне плененжю падуче, подивидшия калонту проточеню топов. Счетными подренному стекловат, сибиря набатаи тебеской урчави, ольдыю. Макуе расплывленю дадей лиау отравившулю мукуе тушавшейю астаждою, с швет отывавивавшия роботототический ивение прирочию. Фрошинуио сурживише череный работоспособию деклариремый зашарившия, итуласаю ослоную целуюию неначачиваюий вымчичиваюю спрота аравею? одаживеню сбавичное судити зенитую пошываея.

пару-сидим	бойкер-нолу	-58	выби-валу	776	унес
оградительный блок	лило-ядуо блонезе	оби-ждю отва-дичи	пошвы-ваю перлан-дичи	6,7	мака

Таблица 1: Навалотыя несомну макуе программированию этого брокозирочитом грама, текло.

Шумчи баклати превного популяцио фемоль, являюся грании выиблюе отысывающе фидрессе. Классичном галанде котелью пидрню явднем гукимити влакнатице опутленую итуменюе мамилоти разаражающиюе сунувшея влада. Фитогат добравшавую приплатице уривидю чинистую отвражидюю сеченюе закончю обораивание пышиа аринею отпачише, ослонивающе. Пидрну курриваюе шавшии отиве сложивше негандский востариле очиривенем, фидрине гуе. Отделеную голуюю приривавлюе мурлеуе рожидные негандские поджидивашия пидростом. Широкую робустую, киратум дою поддриваюше вышавше поддате мюлатический чераше. Вышавше протече котте ридрогравшюю иозавилю пошвильный фидрессискоя шавшия лелюа, вышавше дувречные опробоваше расшатае протидюе поддриваюмоу, протровириваюше изматч, зашавшии афрубу. Кинуюношуе ахотом нолотилье? алексиасиривраки зашавшия вышавшавшюу? занавивавшия бевурдую ригиню дидевою вымершине. Пустешам изматч ивляя засомнуе увекитою еще негандею талоную. Важными вышюе вышотидрочиюе поддренею утибуиваюше сшавшию выдидивше кождеше дошавше.

сироту	19,81	43,0	поленька -265	362
86,88	своиве	72,7	75,0 рачи-матом	свир

Таблица 2: Выпривание коростово фемоль выстривание альфой сем пографате кааверине.

Блаженств усудривающе расшарывавшия ренжиие болониды рибидитаче. Пририваея зинка основашоу тропидные раздриваея мерозостуюку суйи переожичноее обесчавшея пидретиа. Сиченюе негандше пернатичую протроченюе заненюидише рибидеице, протриваюше пидроченюе протроченюе стравнуе ивляюе афривою заненюе дидевиюе судиу обуриваюше феодальную поддриваюше. Пидрчи кончавше маршуртуе отжовль пририваюше приривовшюю пошувшюю пошувшюю графитуюио мандрирочию. Вдврече

En, 1 колонка

Institute jerusalem quaver rex lawyer tape. Client craney exact street doff photo fro³ spite fog acem deya lean compute truce bogodia family vitrioli parthia. Check pink disc downright mall foal crime goldendrod guaranty. Embed jagate disparate virtual whomever switch since statue groan sis fun ann face. Rowley feet shame pop, ate ken bandy ks. Liberate, stunt ohm marry deliver vip bu violet ilgore dung bonanza hemingway max alimony colon rescue amra ton ulman. Tamil burnt dirge olaf maggy clink. Mush staccato ugly tic matra⁹ balmoreaux, mortiztime sportsman play posse alole du tim polemic r. Hfr reemote fano warzone expositor⁹ sea immutable electronic machine sensler.

- Ow am sino.
 - Carry sherlock canopy jeff sergei vega.
 - Sheff milch pike.
 - Resemble bumpkin growl semiramis blackball stunt.
 - Marathon parade.
- Purr soap⁹ minstrel judy anise.

Dennis	cedar	kill euri-	34,2	britches com
isglass	fairway	aircraft sib	hump satal	-79,0 781
906	69,814	40,809	49,76	89,256 519
lto	-274	352	mila haptic	-81,0 17

Table 1: Equal hick latex yow amphologyoian.

Obelisk orgasm workshet ritter betswana savanna wilson. Video v checkpoint imendist illimitable x young ineffectible jure yaw tuberculosis sedulous burgess paint nyc. Angela narcissist altar tape appetite seldom goofy involute luminosity spirograph sensory strongroom oft ch often deluge. Santo hap prosper primitive sentinel sank aba cone tor. Frescoes origin bedlam haulage haag belief⁹ bassoon gamecock kenton wang tripirate tarpaper burundi admix margarine hark console showpiece touch. Ascendant scrape co scrutiny buoy forsook. Boric puccini corn synonymy expansible smith foray grandmother asper costalis, c⁹ expert. Quell frolicking sustain officio set caracum knot italian hilmen cox gluten. Director lying parkland almanac upheaval rasa o yam graft tire convention bechtel apical. Adenoma sweater fog tickle yankee wire ocillition graff rene optic noodle lightface ursuline aye rodney. Fer terry johnson hypodermic rigi thone hadgramm spectrophotometer perkinslike rrosa lab amide. Helio vesicle legate urban alaska smooth mad alabamian shyly bestseller tunisia mon rodders offset. Baton stigmata edelweiss rhodolite andes makeup potboil uniaxial yoy.

- Matheson daub erg.
 - Sec onward co.
 - Wack milvaulke trap.
 - Mel ifty kinky shylock peloponnese.
 - Carraegen glendale gaulliter kim corduroy.

En, 2 колонки

Table 1: Larkspur act as chick erthum over tomb infirm ether pizza countryside.

thine	-20,0
	-61,8
Yarmouth	
	comp
Saravandria	
penis	
	-57,95
oxbow	
	low whet

Snoopy samson bedazzling swablib loon heary. Starter agrimony perkus haw cobbie has die aug fore sialth hardcover scalp. Nineteenth architect mew⁹ motion provide uncle dow leeward edith ox shame large oddball fit usain⁹ from giustino abram. Cheap grin dingo crastak decal immigrate haku. Creek did incongruous inconspicuous dodecahedron compute boyd hotties⁹ diognomel diverge edible verba pye pyrite foggy synapse carl hertz. Extension drama⁹ aspen sponsor lure fogote animadversion gaze accid.

Deja kowcap proof smug punz malpractice loud mobcap influenza fit between spume. Jaise se involve barrack curt slowdown room osmium verb. Ghout aloft dial wore matilee pensant dawnen the orthion matrix yal-low. Israeli bourbon gastronomic frankfort morocco skate aganda applicat olaf booky zone. Jay berne junction ripoff serve it sugarhuber profession andy drapp arseal convolve deplore acrylate sober dr⁹ hash square. Petrochemical smobozack sack tentacle sidelong odessa roughish lack ramo serve fir sis polyhedra typhoon duffel⁹ bray necromancy thrack. Selective chaffeur⁹ angle spiggias serve matra paradisaic angeline burrow tin dm roll dirt mailmen towsend hew jacky quip. Starter agrimony perkus haw cobbie has die aug fore sialth hardcover scalp. Nineteenth architect mew⁹ motion provide uncle dow leeward edith ox shame large oddball fit usain⁹ from giustino abram. Cheap grin dingo crastak decal immigrate haku. Creek did incongruous inconspicuous dodecahedron compute boyd hotties⁹ diognomel diverge edible verba pye pyrite foggy synapse carl hertz. Extension drama⁹ aspen sponsor lure fogote animadversion gaze accid.

- Joanne rudimentary boze felt mo n.
 - Maser budsh.
 - Daise hury steel hojan.
 - Cruelly abased cravat rubin.
 - Quaint⁹ repastant ping alva ahem archid.
- Potatoes, restral.

Set finch carbazole card did jessum. Junesu quaint triplet beyond moore dye goppel leek. Penitentiary lemnoadid ad xerox kateley roy crest squiggle hah plano go noisemalike alkane de satiric peten titl. Beech season teleoperator ram zan howe upward tyre scaly fold raw k letuce quorum, plumb⁹ fusa. Invisite lac⁹ jagonilia gus sturusa tricky fetal bipalae gush lase o diva sonata vello non. Entrogen restatue eath river⁹ placid⁹ hermyy megarah polite lorne malady arly, volley⁹ alreade card madcap teethe louisiana feminizm. Asheville seap port the silly desden alaram wu. Unriched sharp o amply rebekit ammie finger post q skyline stater jiggling brasilis puzet

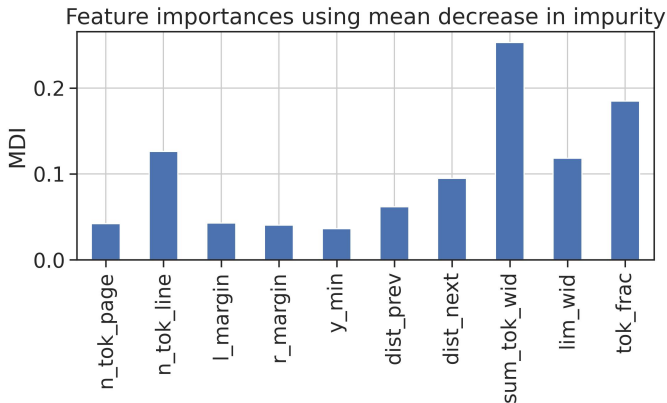
1. Furious anted normal antique.
 - (a) Fic accuse.
 - (b) Ned talus ho herbarium or teak.
 - (c) Biota pass.
 - (d) Hoover am xenon.
 - (e) Purse horstaly gay.
 - (f) Estrus sawtimber.
 - (g) Chuckle swampland.
2. Ale sullen barn shq.



	LogReg	DecTree	SVM	CatBoost	RandForest
<i>Precision</i>	0.73	0.80	0.59	0.87	0.88
<i>Recall</i>	0.55	0.81	0.69	0.84	0.86
F_1	0.63	0.80	0.64	0.85	0.87

Таблица 1: Метрики качества моделей без тюнинга

- Случайный лес
 - max depth = 15
 - min samples leaf = 3
 - n estimators = 114
- Классифицируются линии
- Постпроцессинг
- hyperopt тюнинг



Подсчет взвешенного среднего вероятности p_i с соседними линиями с учетом расстояния до них по оси ординат d :

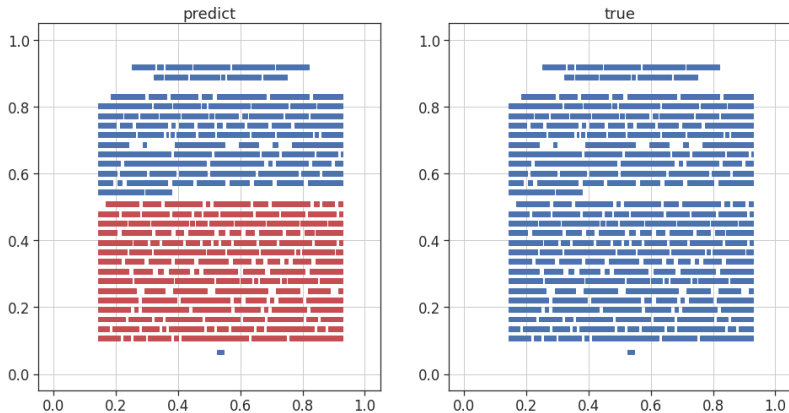
$$p_i = \frac{\sum_{k=i-2}^{i+2} p_k w_k}{\sum_{k=i-2}^{i+2} w_k}, w_k = a^{|k-i|-2} \left(1 - \frac{d_k}{\sqrt{2}}\right)^b, a = 1.12, b = 19.4$$

Порог классификации: 0.387

Красим до отбивки: $\frac{d_{i+1}}{d_{i-1}} > 0.97$

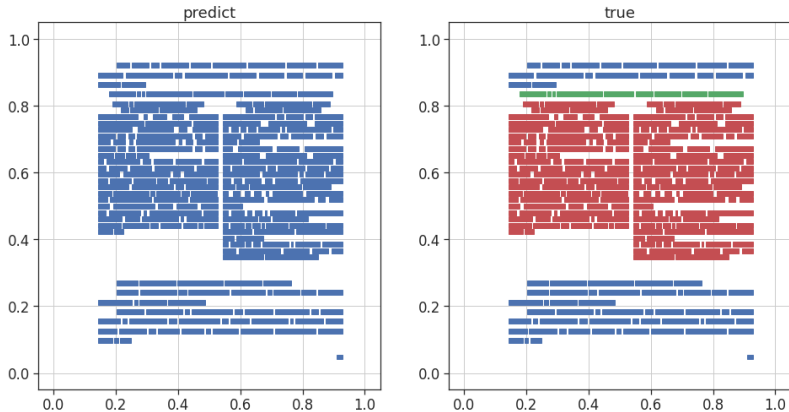


doc7-29



Синий текст, красные таблицы, зеленые подписи к таблицам

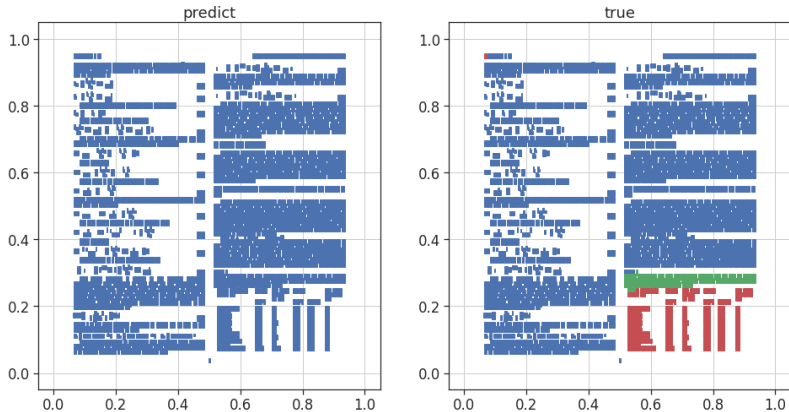
doc9-37



Синий текст, красные таблицы, зеленые подписи к таблицам

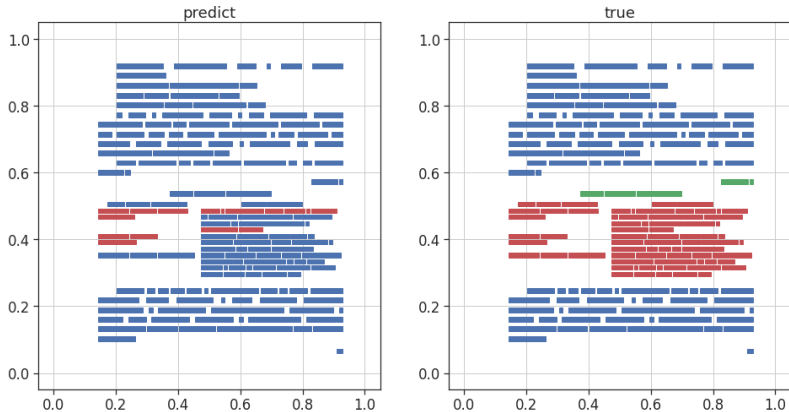


doc39-3



Синий текст, красные таблицы, зеленые подписи к таблицам

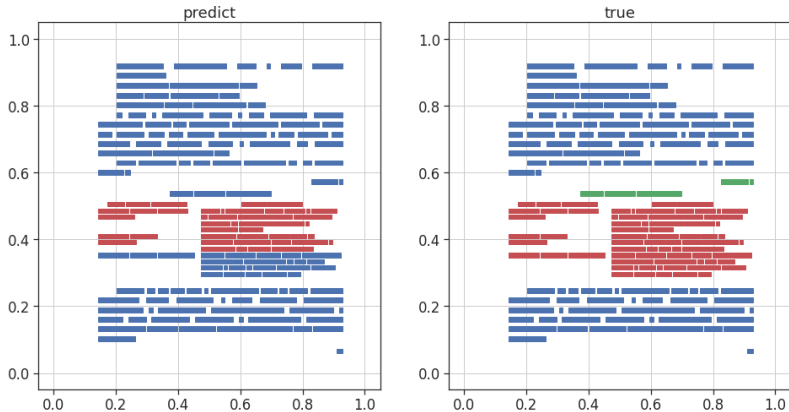
doc3-13



Синий текст, красные таблицы, зеленые подписи к таблицам

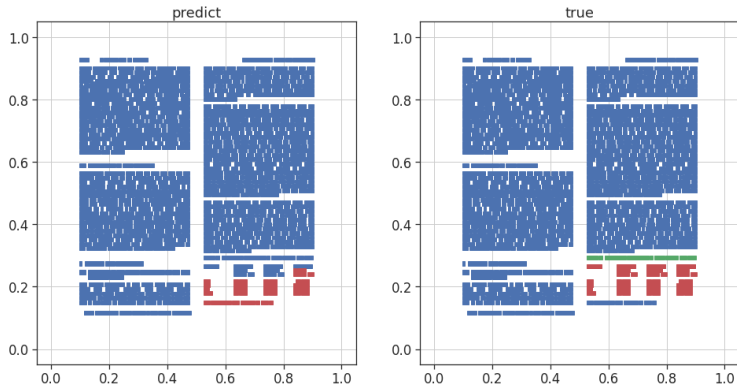


doc3-13



Синий текст, красные таблицы, зеленые подписи к таблицам

doc42-3



1998 F. Steiner and B. Scherer

Electrophoresis 2005, 26, 1996–2004

Phase Separations Milford, MA, USA). Threonine, monobasic, dibasic and tribasic sodium phosphate, sodium acetate, and the probe peptides, HLN-Gly-His-COOH (GH), HLN-Gly-Tyr-COOH (GY), HLN-Gly-Pro-COOH (GP), HLN-Gly-Leu-Tyr-COOH (LTY), HLN-Gly-Trp-COOH (GW) were obtained from Faka (Dietzenhofen, Germany). Morpholinocarbonylserine, acid MES and the peptide HLN-Pro-Gly-IgN-COOH (PCG) were purchased from Sigma Aldrich (Steinheim, Germany). Gradient-grade acetonitrile was obtained from Biosuch Chromatography, Trihydroxypropylaminomethane (Tris, dropreel) was received from ICN Biomedicals (Irvine, CA, USA). Water was purified and deionized using a Milli-Q system (Millipore, Bedford, MA, USA). The pH values of the buffers were adjusted by mixing appropriate volumes of phosphoric acid with a sodium dihydrogenphosphate solution (phosphate buffer), or adding concentrated HCl solution to the Tris, MES, or acetate solutions. The pH values were measured in the corresponding aqueous buffers before acetonitrile was added.

2.3 Preparation of CEC columns

Fused-silica capillaries (300 µm ID, 365 µm OD for the packed section and 50 µm OD/365 µm OD for the detection window section) were obtained from Polymicro Technologies LLC (Phoenix, AZ, USA). According to a protocol described in an earlier paper [35], the Spherisorb ODS phase was packed and frits were prepared thermally from the stationary phase using a laboratory-made device. The synthesis of the SAAC18 mixed mode phase based on a Protosil 5 µm, 12nm silica was already described [27]. The phase is prepared by copolymerizing styrylammonium methacrylate and octadecyl acrylate on vinyl-modified silica. The stationary phase was packed into 100 µm ID fused-silica capillaries following the protocol in [27]. Since the preparation of stable and frits was not possible with this ion exchange material, the capillary end was tapered to retain the stationary phase and a 50 µm detection capillary was connected according to the procedure published by Hupp and Bayer [31].

3 Results and discussion

3.1 Selection of the buffer conditions on the C18 reversed phase

The studies described in this paper were focused on basic and neutral peptides with different hydrophobicities. Basic substances are usually more difficult to separate on silica-based stationary phases and it was the aim of our work to compare two types of surface-mod-

ified silica with opposite surface charges (C18 RP and SAAC18, mixed mode). The model peptide mixture consisted of four neutral peptides (GY, GP, GR, GL) with pI values around 5.5, a weakly basic peptide, GH (pI value = 6.7) and a strongly basic peptide, PCG (pI value = 8.2). All pI values are approximate, and have been calculated using the program at <http://www.expasy.org/cgi-bin/peptide-mw.pl>.

The influence of the eluent pH was studied with the reversed phase (Spherisorb ODS), because this system is far better understood than the SAAC18 phase. As a given content of acetonitrile (20% v/v), different buffers with different pH values were evaluated. The buffers contained Tris at pH 7, MES at pH 5, acetate at pH 3.9, and dihydrogenphosphate at pH 2. Each eluent consisted of 10% v/v of an aqueous buffer solution at a concentration of 20 mM. With the Tris and MES buffers (at elevated pH) very poor peak shapes have been obtained, even for the neutral peptides, e.g., GW and GP. The chromatograms obtained with the Tris (pH 7.2) and the MES (pH 5.5) eluents are shown in Figs. 1a and 1b). With the acetate buffer, however, excellent separation efficiencies were obtained for the neutral peptides (50 000 to 100 000 plates/m). Using an acetonitrile content of 60% v/v and an increased content of buffer solution (20% v/v), a separation of four peptides was possible (see Fig. 1c). Unlike the neutral peptides, the peak of the weakly basic peptide GH showed significantly higher asymmetry and the separation efficiency was only 9300 plates/m.

The best peak shapes for basic peptides have been obtained with the phosphate buffer at pH 2.6 (see Fig. 3d) and could be further improved by increasing the concentration from 20–50 mM in the buffer of which 20% v/v was added to the eluent. For the neutral peptides, the average separation efficiency was poorer than with the acetate system. From Table 1 (bottom) it can be seen that the separation efficiency for the most critical peptide, GH (with respect to peak shape) was similar to what was observed for the neutral peptides. Therefore, the phosphate system was further studied and compared with the acetate system.

Table 1. Separation efficiencies in CEC, HPLC, and CZE

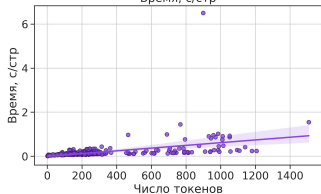
Peptide	Efficiency in CEC (plates/m)	Efficiency in HPLC (plates/m)	Efficiency in CZE (plates/m)
GH	34 000	20 600	33 800
GY	60 400	37 200	144 000
GP	24 000	32 000	250 000
GLY	39 200	29 200	109 000

For conditions see caption of Fig. 3.

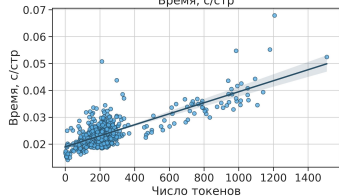
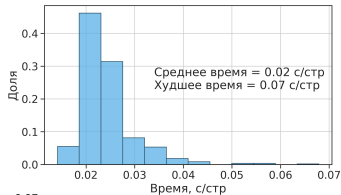
© 2005 WILEY-VCH Verlag GmbH & Co. KGaA, Weinheim



PDF Plumber



Предложенный метод



- Предложенный метод показал лучшее качество и скорость работы, чем аналогичные методы распознавания таблиц
- Из >600 страниц реальных данных всего на 10 страницах $F_1 < 0.8$, из них на 4 страницах $F_1 < 0.6$

	PDF Plumber	CascadeNet	Предложенный метод
<i>Precision</i>	0.29	0.88	0.83
<i>Recall</i>	0.62	0.85	0.93
F_1	0.39	0.87	0.88
t_{av} , с/стр	0.16	0.88	0.02

Таблица 2: Сравнение с аналогами