# EXTRACTING TABLE DATA FROM PDF

### A JOURNEY INTO THE UNKNOWN

by Don Vito Perleone

German Perl Workshop 2016 Nürnberg

## **BASED ON A TRUE STORY**

# **ABOUT ME**

- 20 years of commercial Perl experience
- Enjoy messing about with data
- Hannover.pm

# DR. STRANGELOVE — WHY?

### **REAL-WORLD EXAMPLE**

- Cafeteria at \$work, run by external service provider
- Weekly menu available as PDF

Speisep	lan I	KW	05

	Montag Dienstag		Mittwoch	Donnerstag	Freitag	
	1. Februar 2016	2. Februar 2016	3. Februar 2016	4. Februar 2016	5. Februar 2016	
Suppe 0,85 €	Tomatencremesuppe	Gemüsebrühe <sup>28</sup>	Kartoffelsuppe <sup>28</sup>	Zwiebelsuppe	Erbsencremsuppe <sup>26,28</sup>	
Haupt- gerichte	Koreanisches Feuerfleisch vom Schwein mit Mangold und Wildreis 3,60 €	Riesencurrywurst <sup>2,3,8</sup> mit zweierlei Saucen <sup>1,2</sup> und Pommes Frites 3,60 €	Gebratene Hähnchenbrust in Rosmarin-Tomatensauce Nussreis <sup>27</sup> und Grillgemüse <sup>28</sup>	Gebackener Leberkäse <sup>20w</sup> mit Bohnengemüse <sup>26</sup> und Kartoffelstampf 3,60 €	Schlemmerfilet á la Bordelaise <sup>4,20</sup> mit Dillschaum <sup>20,26</sup> Tomatensalat und Kartoffeln 3,60 €	
genome	Putengulasch mit Zucchini und Champignon <sup>20w,26,28</sup> dazu Nudeln <sup>20w,23</sup> 3,60 €	Kutterfrikadelle <sup>4,21</sup> Gurkensalat <sup>2,3</sup> und Kartoffelpüree 3,60 €	Hausgemachtes Schnitzel <sup>20w</sup> "Mailänder Art" mit Spaghetti <sup>20w</sup> und Tomatensauce 3,60 €	mit Rotweinjus, Rosenkohl und auf	Tanduri Masala Putenbraten auf Gemüse-Cous-Cous mit Joghurtdip <sup>26</sup> 3,60 €	
Vegeta- risch	Fit Menü  Gratinierte Maultaschen <sup>20(W),23,28</sup> auf Tomatenspinat	Fit Menü  Frühlingsrolle <sup>2,4,20,25,30</sup> mit Süß-Saurem Dip und Basmatireis	Fit Menü  Gemüselasagne <sup>20w,26,28</sup> mit  Basilikum - Tomatensauce	Fit Menü Maisplinsen <sup>20w,23,26</sup> an gegrilltem Gemüse mit Limettensauce <sup>26</sup>	Fit Menü  Rahmige Kartoffelsuppe <sup>26,28</sup> mit bunter Gemüseeinlage und Baguette <sup>20,23,25,26,28,29,30,32</sup>	
Salate	3,30 € 3					
Beilagen ₀,70 €	Wildreis Nudeln Tagesgemüse <sup>28</sup>	Pommes Frites Kartoffelpüree Basmatireis Tagesgemüse <sup>28</sup>	Nussreis Spaghetti Tagesgemüse <sup>28</sup>	Kartoffelstampf Spätzle Tagesgemüse <sup>28</sup>	Kartoffelsalat Reis Tagesgemüse <sup>28</sup>	
Dessert Groß 0,80 € Klein 0,50 €					<u> </u>	

Klein 0,50 €

Zusatzstoffe sind mit Zahlen gekennzeichnet siehe Aushang

Bei Fragen zu Dussmann können Sie auf www.dussmann.com klicken!

Speiseplanänderungen vorbehalten

Speiseplan KW 06						
	Montag	Dienstag	Mittwoch	Donnerstag	Freitag	
	8. Februar 2016	9. Februar 2016	10. Februar 2016	11. Februar 2016	12. Februar 2016	
Suppe 0,85 €	Hühnerbrühe mit Einlage	Käse - Lauchsuppe	Minestrone	Karottencremesuppe mit Ingwer <sup>26,28</sup>	Pfifferlingcremesuppe	
Haupt- gerichte	Gefüllte Paprikaschote mit Schweinemett, Reis und Tomatensauce 3,60 €	Crispy Chicken Burger mit Asiasauce und Pommes Frites <sup>20,23,26</sup> 3,60 €	Kohlrouladen mit Hackfleischfüllung, <sup>20w,26,29</sup> Kohlgemüse und Salzkartoffeln 3,60 €	Griechische Bratkartoffeln mit Hähnchenbruststreifen Zucchini und Hirtenkäse <sup>26</sup> 3,60 €	Grünkohl <sup>20H,28,29</sup> mit Kassler, Mettende <sup>1,2,3</sup> und Kartoffeln 3,80 €	
genente	Ente Süß-Sauer <sup>1,2,3,9,20w,25,28,29</sup> mit gebratenen Nudeln 3,60 € Fit Menü	Gedünsteter Tilapia <sup>21</sup> "Lingurische Art" mit Oliven und Tomaten <sup>28</sup> dazu Pasta <sup>20w,23</sup> 3,60 € Fit Menü	Grüne Thaisuppe <sup>20</sup> mit Hähnchenbrust und Reisnudeln 3,60 € Fit Menü	Highlight der Woche  Gebratene Forelle auf Gemüse- Kartoffel-Bett mit Tomaten - Mozarellasauce  5,20 €  Fit Menü	Deftiger Bohneneintopf mit Kassler und Baguette <sup>20w,26,28</sup> 3,60 €	
Vegeta- risch	Blumenkohl "Indisch" mit Nüssen, Curry, Rosinen <sup>24,26,27</sup>	Milchreis <sup>26</sup> mit Heißen Kirschen und Zimt Zucker	Penne <sup>20</sup> Napoli und Grana Padano <sup>23,26</sup>	Ofenkartoffel mit zweierlei Quark <sup>26</sup> dazu Salatbeilage	Gekochte Eier <sup>23</sup> in Dijonsenfrahm <sup>20w,28,29</sup> mit Salzkartoffeln	

Salate

Bei der Zubereitung der Speisen verwenden wir Allergie / Lebensmittelunverträglichkeit auslösende Stoffe als Zutaten, diese sind im Speiseplan gekennzeichnet. Trotz größtmöglicher Sorgfalt können wir jedoch nicht ausschließen, dass technisch unvermeidbare Einträge von Allergie / Lebensmittelunverträglichkeit auslösenden Stoffen und ggf. weitere mögliche Einträge auch in andere Gerichte gelangen. Weitere Auskünfte dazu erteilt Ihnen gerne unser Küchenpersonal.

3,30 €

Beilagen ₀,70 €
0,70 €

Reis gebratene Nudeln Basmatireis Tagesgemüse<sup>28</sup>

und Basmatireis

3,30 €

Pommes Frites Pasta Tagesgemüse<sup>28</sup>

3,30 €

Kartoffeln Reisnudeln Penne Tagesgemüse<sup>28</sup>

Griechische Bratkartoffeln Ofenkartoffel Tagesgemüse<sup>28</sup>

3,30 €

Kartoffeln Reis Tagesgemüse<sup>28</sup>

und Rote Beete Salat9

3,30 €

Dessert Groß 0,80 € Klein 0,50 €

Zusatzstoffe sind mit Zahlen gekennzeichnet siehe Aushang Bei Fragen zu Dussmann können Sie auf www.dussmann.com klicken!

Speiseplanänderungen vorbehalten

### Speiseplan KW 05

	Montag	Dienstag	Mittwoch
	1. Februar 2016	2. Februar 2016	3. Februar 2016
Suppe 0,85 €	Tomatencremesuppe	Gemüsebrühe <sup>28</sup>	Kartoffelsuppe <sup>28</sup>
Haupt- gerichte	Koreanisches Feuerfleisch vom Schwein mit Mangold und Wildreis	Riesencurrywurst <sup>2,3,8</sup> mit zweierlei Saucen <sup>1,2</sup> und Pommes Frites 3,60 €	Gebratene Hähnchenbrust in Rosmarin-Tomatensauce Nussreis <sup>27</sup> und Grillgemüse <sup>28</sup>
	Putengulasch mit Zucchini und Champignon <sup>20w,26,28</sup> dazu Nudeln <sup>20w,23</sup> 3,60 €	Kutterfrikadelle <sup>4,21</sup> Gurkensalat <sup>2,3</sup> und Kartoffelpüree 3,60 €	Hausgemachtes Schnitzel <sup>20w</sup> "Mailänder Art" mit Spaghetti <sup>20w</sup> und Tomatensauce 3,60 €
	Fit Menü	Fit Menü	Fit Menü
Vegeta-	Gratinierte Maultaschen	Frühlingsrolle <sup>2,4,20,25,30</sup>	Gemüselasagne <sup>20w,26,28</sup> mit

### Speiseplan KW 06

	Montag	Dienstag	Mittwoch
	8. Februar 2016	9. Februar 2016	10. Februar 2016
Suppe 0,85 €	Hühnerbrühe mit Einlage	Käse - Lauchsuppe	Minestrone
Haupt- gerichte	Gefüllte Paprikaschote mit Schweinemett, Reis und Tomatensauce 3,60 €	Crispy Chicken Burger mit Asiasauce und Pommes Frites <sup>20,23,26</sup> 3,60 €	Kohlrouladen mit Hackfleischfüllung, <sup>20w,26,29</sup> Kohlgemüse und Salzkartoffeln 3,60 €
genome	Ente Süß-Sauer <sup>1,2,3,9,20w,25,28,29</sup> mit gebratenen Nudeln 3,60 €	Gedünsteter Tilapia <sup>21</sup> "Lingurische Art" mit Oliven und Tomaten <sup>28</sup> dazu Pasta <sup>20w,23</sup> 3,60 €	Grüne Thaisuppe <sup>20</sup> mit Hähnchenbrust und Reisnudeln 3,60 €
	Fit Menü	Fit Menü	Fit Menü
Vegeta-	Blumenkohl "Indisch" mit Nüssen, Curry,	Milchreis <sup>26</sup> mit Heißen	Penne <sup>20</sup> Napoli und Grana

### FOOD RATING APPLICATION "GOURMETER"

- Made by our software engineering trainee
- Needs weekly manual input of the menus
- Luckily not by me
- But I do enjoy scripting things
- and making them run on their own

### **AUTOMATION BECKONS**

- Not as easy as it looks
- Naive approach doesn't work
- ... at all
- Dig deeper
- Much deeper

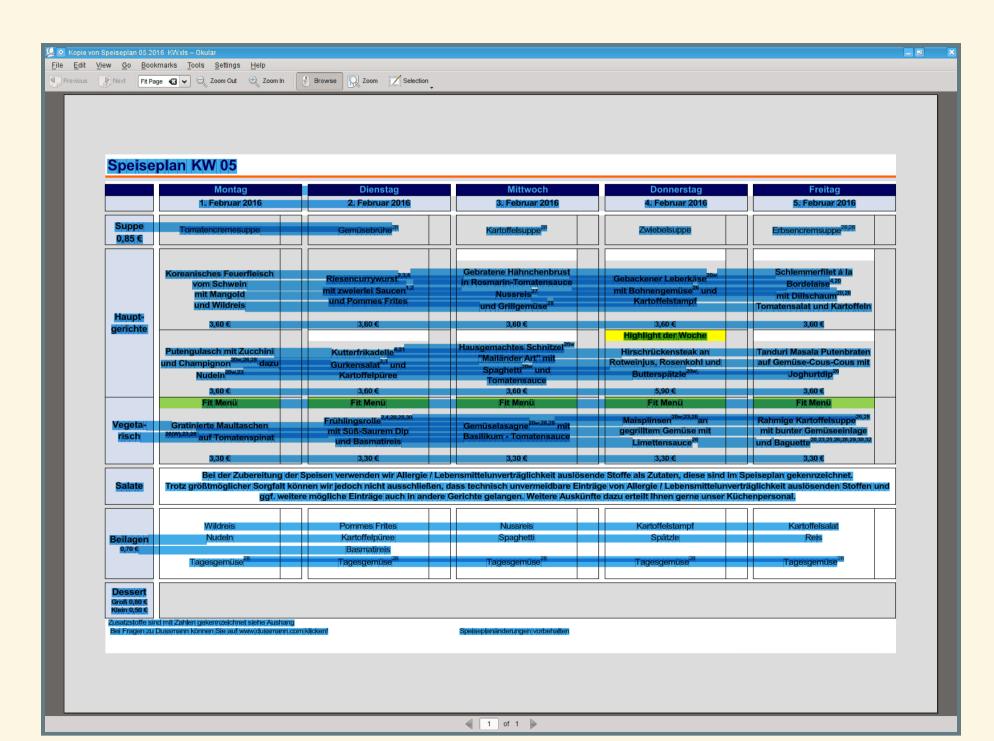
### WHY DID I DO THIS?

- London Perl Workshop last December
- Attended talk "Perl and PDF" by Prabhakar Somu
- http://act.yapc.eu/lpw2015/talk/6423
- https://youtu.be/5\_yFwl5zN8w
- Motivated me to submit, research and prepare talk

# FIRST BLOOD - NAIVE APPROACH

### **COPYING THE TEXT**

- It's right there, I can see it
- I can also select it



### **COPY TEXT**

- It's right there, I can see it
- I can also select it
- And copy it
- But the result is
- not quite what you expect

```
Speiseplan KW 05
 2
3
     Montag
      Dienstag
 4
     1. Februar 2016
 5
      2. Februar 2016
 6
7
     Suppe
       28
 8
     Tomatencremesuppe
 9
      Gemüsebrühe
10
     0,85€
11
     Mittwoch
12
     Februar 2016
     Kartoffelsuppe
13
14
     28
15
     Donnerstag
     4. Februar 2016
16
17
     Zwiebelsuppe
     Freitag
18
19
     Februar 2016
20
     Erbsencremsuppe
21
     26,28
22
     Koreanisches Feuerfleisch
23
      2,3,8
24
      Gehratene Hähnchenhrust
```

A. T	ocbi a cerie i farificii ci bi a se
25	20w
26	Schlemmerfilet á la
27	Riesencurrywurst Gebackener Leberkäse 4,20
28	vom Schwein

### **COMMAND LINE TEXT EXTRACTION**

- Several tools
  - ps2txt (ghostscript)
  - pstotxt (also ghostscript)
  - pdftotext (xpdf, poppler)
- Output ranges from total crap
- to all there, but wrong order

```
1
2
3 \Delta \Theta \Lambda \Lambda \Xi \Pi \Sigma \Upsilon \Phi \Psi \Omega
4
5 fffifl\Xi ffifflij`'\Xi j\Psi ^\Xi \Theta \Xi fl*_\Xi ij`'*fi,\Psi \Delta `'ss
6
7 \Theta fflAE\Psi OEi_AEfl\Xi ij
8
9 O/i\Xi j\Xi ffl`\Theta flflss\Theta fljae!\Sigma "\Sigma \Upsilon ,iae\Psi #ss
10
11 \Theta fflAE\Psi &fli__o/\Xi ,.j\Xi !\Upsilon
12
13 &\Xi 'ffi`/\Xi ffl\Xi fl\Psi O\Xi '\Xi fl/)j\Xi !\Pi ss\Psi ,iae\Psi 1fi'ffl\X
```

```
1 12334
2 56789\0mega
3 fffifl4ffifflij`'4j9<9F>424fl<96>4ij`' <9A>fi91`'4iffl iae9oeffifflfiAE 2fflAE
4 i4j4ffl`2flfl/2fljae !6''67
5 iae9#4i4fl4i91ffi2`4ffl $6!
6 2fflAE9%fi4i9<9F>fliae4i
7 &4'flffiae4ffl49()'ffl`'4ffl'fl2jae
8 iffl9fijffifliffl*+fiffiae4ffljffi2`4 ,2jjfl4ij !
9 2fflAE9&fli4.j4 !7
10 &4'ffi\'/4ffl4fl904'4fl/)i4 !59
11 iae91fi'ffl4ffl4.j4 !2 92fflAE9
12 ffffiflaefi<96><96>4jaeffi3<96>
13 1`'44fl<96>i4ae939ffi9 1fiflAE4ffiij4 46!5
14 iae95ii`'ffi2 !56!29
15 +fiffiae4ffliffiffiae92fflAE9ffffiflaefi<96><96>4ffl
16 ''6259\Omega ''6259\Omega ''6259\Omega ''6259\Omega ''6259\Omega (i'i'ae9AE4fl
17 %2ae4ffl2ffij`'9iae962``'iffli9
18 2fflAE97'ffi3ifflfiffl !56!26!7 9AEffi#29
19 ,2AE4ffl !56!''
20 ff2aeae4fl<96>fli/ffiAE44 46!$9
21 &2fl/4ffliffiffiae !6'' 92fflAE9
22 ffffiflaefi<96><96>43.fl44
23 (ffi2j4ffi`'ae4j91`'ffliae#4 !59
24 8oeffii)fflAE4fl99flae89iae9 13ffi'4aeaei !5 92fflAE9
```

```
Speiseplan KW 05
                                           35 Koreanisches Feuerfleisch
                                           36 vom Schwein
    Montag
                                           37 mit Mangold
                                           38 und Wildreis
    Suppe
 5
    0,85 €
                                           39
6
                                           40 Hauptgerichte
    Dienstag
                                           41
8
                                           42 Riesencurrywurst2,3,8
9
    Mittwoch
                                           43 mit zweierlei Saucen1,2
10
                                           44 und Pommes Frites
11
                                           45
    Donnerstag
12
                                           46 Gebratene Hähnchenbrust
13
                                           47 in Rosmarin-Tomatensauce
    Freitag
14
                                           48 Nussreis27
15
    1. Februar 2016
                                           49 und Grillgemüse28
16
                                           50
17
    2. Februar 2016
                                           51 3,60 €
18
                                           52
19
    3. Februar 2016
                                           53 3,60 €
20
                                           54
21
                                           55 3,60 €
    4. Februar 2016
22
                                           56
23
    5. Februar 2016
                                           57 Zwiebelsuppe
24
                                           58
                                           59 Erbsencremsuppe
    Tomatencremesuppe
```

### **QUICK SIDEBAR**

### HOW I BUILT THE PREVIOUS TWO-COLUMN SLIDE

```
$ paste -d=
<( cat -n menu.txt | sed -n 1,25p )
<( cat -n menu.txt | sed -n 35,+24p )
| column -s= -t | expand</pre>
```

### PDF ...

- Short for Portable **Document** Format
- Concerned with Looks and Layout
- Structure?
- What structure?

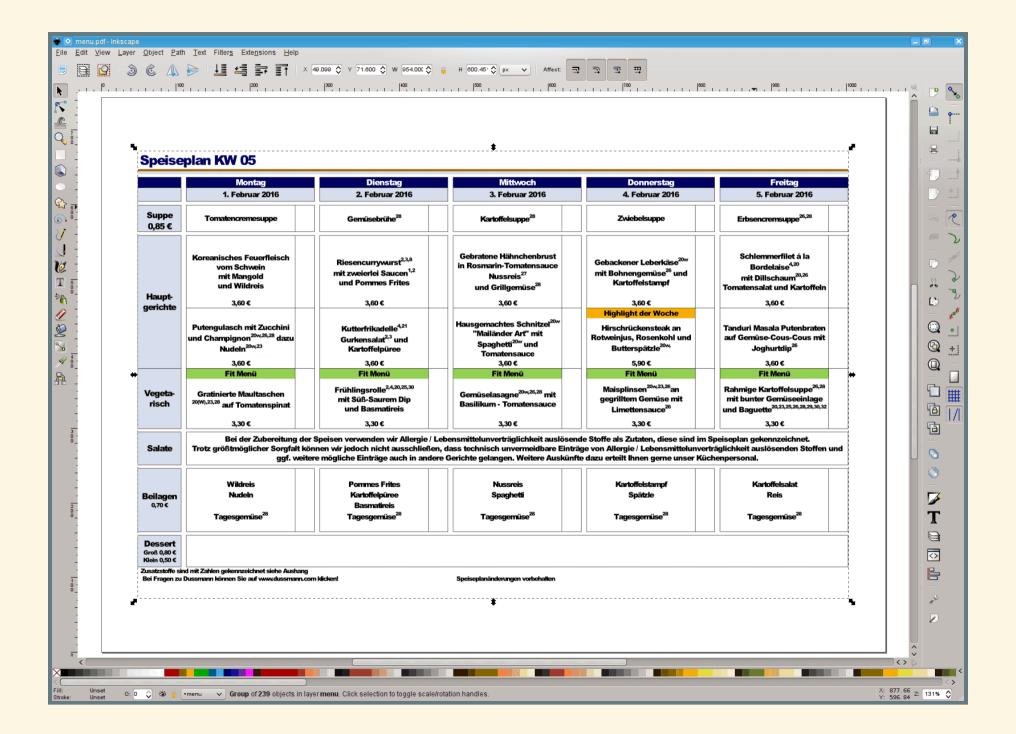
# SAVING PRIVATE RYAN — SNEAKY APPROACH

### WE'RE GONNA NEED A BIGGER CALIBER!

- SVG
- Libre Office
- Last Resort: OCR

### SVG

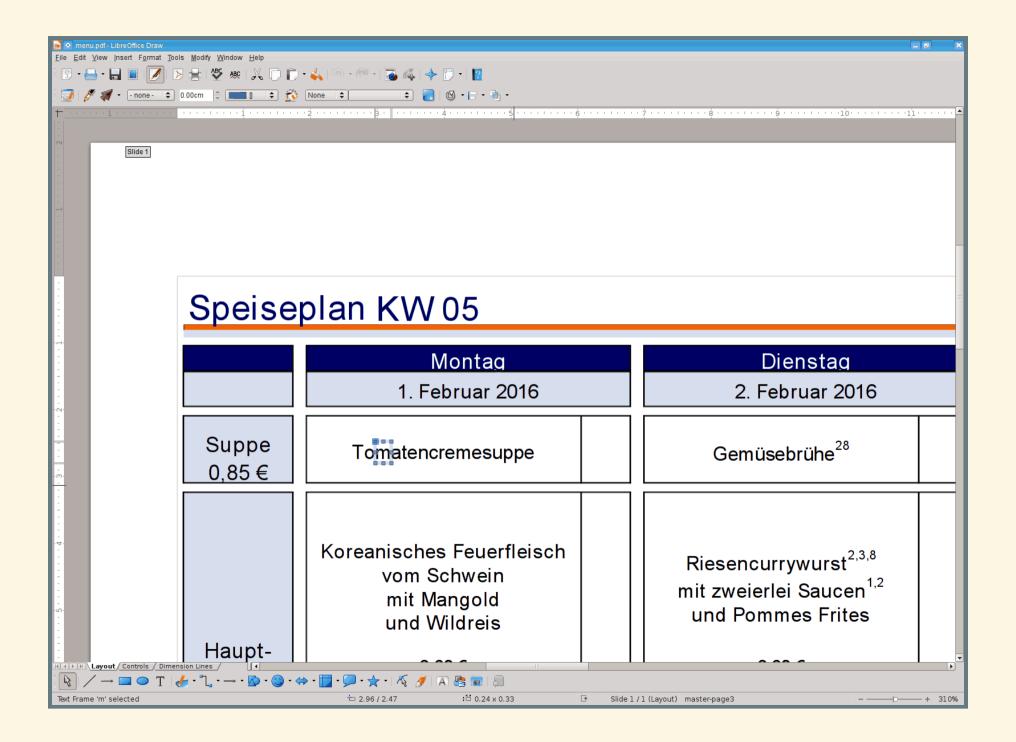
- Text based vector format
- Inkscape
- ImageMagick
- pdftocairo (poppler)
- Grouping not useful



	Montag	Dienstag	Mittwoch	Do	
<u> </u>	1. Februar 2016	2. Februar 2016	3. Februar 2016	4. F	
Suppe 0,85 €	Tomatencremesuppe	Gemüsebrühe <sup>28</sup>	Kartoffelsuppe <sup>28</sup>	Zwieb	
Haupt-	Koreanisches Feuerfleisch vom Schwein mit Mangold und Wildreis	Riesencurrywurst <sup>2,3,8</sup> mit zweierlei Saucen <sup>1,2</sup> und Pommes Frites	Gebratene Hähnchenbrust in Rosmarin-Tomatensauce Nussreis <sup>27</sup> und Grillgemüse <sup>28</sup>	Gebackener mit Bohnen Kartoff	
gerichte	3,60 €	3,60 €	3,60 €	3,	
	Putengulasch mit Zucchini und Champignon <sup>20w,26,28</sup> dazu Nudeln <sup>20w,23</sup>	Kutterfrikadelle <sup>4,21</sup> Gurkensalat <sup>2,3</sup> und Kartoffelpüree	Hausgemachtes Schnitzel <sup>20M</sup> "Mailänder Art" mit Spagnetti <sup>20M</sup> und	Highlight Hirschrüc Rotweinjus, Butters	
	3,60 €	3,60 €	3,60 €	5,	
	Fit Menü	Fit Menü	Fit Menü	Fit	
Vegeta- risch	Gratinierte Maultaschen <sup>20(W),23,28</sup> auf Tomatenspinat	Frühlingsrolle <sup>2,4,20,25,30</sup> mit Süß-Saurem Dip und Basmatireis	Gemüselasagne <sup>20,4,26,28</sup> mit Basilikum - Tomatensauce	Maisplins gegrilltem Limette	
	3,30 €	3,30 €	3,30 €	3	

### LIBRE OFFICE

- Can open PDF files
- Imported as drawing
- Each letter a separate element
- Not suited for table extraction



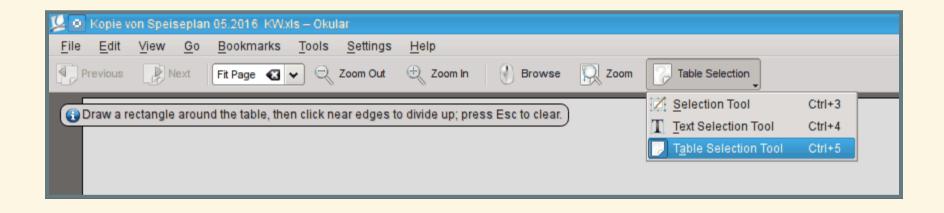
### OCR

- PDF is vector based
- Endless zooming
- PNG with huge, well-defined letters
- Cut up into single cells
- We should OCR those easily (I hope)
- A path not taken

# APOCALYPSE NOW — INTERACTIVE

### **OKULAR**

- KDE document viewer
- Table selection tool



### **OKULAR**

- KDE document viewer
- Table selection tool
- Recognizes row boundaries

	montag	Dictistay	MILLWOCH	Dominiciata	r remay
	1. Februar 2016	2. Februar 2016	3. Februar 2016	4. Februar 2016	5. Februar 2016
	Tomatencremesuppe	Gemüsebrühe <sup>28</sup>	Kartoffelsuppe <sup>28</sup>	Zwiebelsuppe	Erbsencremsuppe <sup>26,28</sup>
	Koreanisches Feuerfleisch vom Schwein mit Mangold und Wildreis	Riesencurrywurst <sup>2,3,8</sup> mit zweierlei Saucen <sup>1,2</sup> und Pommes Frites	Gebratene Hähnchenbrust in Rosmarin-Tomatensauce Nussreis <sup>27</sup> und Grillgemüse <sup>28</sup>	Gebackener Leberkäse <sup>20w</sup> mit Bohnengemüse <sup>26</sup> und Kartoffelstampf	Schlemmerfilet á la Bordelaise <sup>4,20</sup> mit Dillschaum <sup>20,26</sup> Tomatensalat und Kartoffeln
_	3,60 €	3,60 €	3,60 €	3,60 €	3,60 €
ľ				Highlight der Woche	
	Putengulasch mit Zucchini und Champignon <sup>20w,26,28</sup> dazu Nudeln <sup>20w,23</sup>	Kutterfrikadelle <sup>4,21</sup> Gurkensalat <sup>2,3</sup> und Kartoffelpüree	Hausgemachtes Schnitzel <sup>20w</sup> "Mailänder Art" mit Spaghetti <sup>20w</sup> und Tomatensauce	Hirschrückensteak an Rotweinjus, Rosenkohl und Butterspätzle <sup>20w,</sup>	Tanduri Masala Putenbraten auf Gemüse-Cous-Cous mit Joghurtdip <sup>26</sup>
	3,60 €	3,60 €	3,60 €	5,90 €	3,60 €
	Fit Menü	Fit Menü	Fit Menü	Fit Menü	Fit Menü
1-	Gratinierte Maultaschen <sup>20(W),23,28</sup> auf Tomatenspinat	Frühlingsrolle <sup>2,4,20,25,30</sup> mit Süß-Saurem Dip und Basmatireis	Gemüselasagne <sup>20w,26,28</sup> mit Basilikum - Tomatensauce	Maisplinsen <sup>20w,23,26</sup> an gegrilltem Gemüse mit Limettensauce <sup>26</sup>	Rahmige Kartoffelsuppe <sup>26,28</sup> mit bunter Gemüseeinlage und Baguette <sup>20,23,25,26,28,29,30,32</sup>
	3,30 €	3,30 €	3,30 €	3,30 €	3,30 €
	Bei der Zubereitung der S	Speisen verwenden wir Allergie / L	ebensmittelunverträglichkeit auslöse	nde Stoffe als Zutaten, diese sind im	Speiseplan gekennzeichnet.
	Bei der Zubereitung der S	Speisen verwenden wir Allergie / L	ebensmittelunverträglichkeit auslöse	nde Stoffe als Zutaten, diese sind im	Speiseplan gekennzeichnet.

### **OKULAR**

- KDE document viewer
- Table selection tool
- Recognizes row boundaries
- Needs manual indication of column boundaries
- Does the heavy lifting
- Copy/paste as table data

Tomatencremesuppe	28 Gemüsebrühe	Kartoffelsuppe 28	Zwiebelsuppe	Erbsencremsuppe 26,28
Koreanisches Feuerfleisch vom Schwein mit Mangold und Wildreis	2,3,8 Riesencurrywurst 1,2 mit zweierlei Saucen und Pommes Frites	Gebratene Hähnchenbrust in Rosmarin-Tomatensauce 27 Nussreis 28 und Grillgemüse	20w Gebackener Leberkäse 26 mit Bohnengemüse und Kartoffelstampf	Schlemmerfilet á la 4,20 Bordelaise 20,26 mit DillschaumTomatensalat und Kartoffeln
3,60 €	3,60 €	3,60 €	3,60 €	3,60 €
			Highlight der Woche	
Putengulasch mit Zucchini 20w,26,28und Champignon dazu 20w,23 Nudeln	4,21 Kutterfrikadelle 2,3 Gurkensalat und Kartoffelpüree	20w Hausgemachtes Schnitzel "Mailänder Art" mit 20w Spaghetti und Tomatensauce	Hirschrückensteak an Rotweinjus, Rosenkohl und 20w, Butterspätzle	Tanduri Masala Putenbraten auf Gemüse-Cous-Cous mit 26 Joghurtdip
3,60 €	3,60 €	3,60 €	5,90 €	3,60 €
Fit Menü	Fit Menü	Fit Menü	Fit Menü	Fit Menü
Gratinierte Maultaschen 20(W),23,28 auf Tomatenspinat	2,4,20,25,30 Frühlingsrolle mit Süß-Saurem Dip und Basmatireis	20w,26,28 Gemüselasagne mit Basilikum - Tomatensauce	_	26,28 Rahmige Kartoffelsuppemit bunter Gemüseeinlage und Baguette20,23,25,26,28,29,30,32
3,30 €	3,30 €	3,30 €	3,30 €	3,30 €

A	В	C	D D	F
Tomatencremesuppe	28 Gemüsebrühe	Kartoffelsuppe 28	Zwiebelsuppe	Erbsencremsuppe 26,28
2 Koreanisches Feuerfleisch vom Schwein mit Mangold und Wildreis	2,3,8 Riesencurrywurst 1,2 mit zweierlei Saucen und Pommes Frites	Gebratene Hähnchenbrust in Rosmarin-Tomatensauce 27 Nussreis 28 und Grillgemüse	20w Gebackener Leberkäse 26 mit Bohnengemüse und Kartoffelstampf	Schlemmerfilet á la 4,20 Bordelaise 20,26 mit DillschaumTomatensalat und Kartoffeln
3,60 €	3,60 €	3,60 €	3,60 €	3,60 €
			Highlight der Woche	
Putengulasch mit Zucchini 20w,26,28und Champignon dazu 20w,23 Nudeln	Gurkensalat und Kartoffelpüree	20w Hausgemachtes Schnitzel "Mailänder Art" mit 20w Spaghetti und Tomatensauce	Hirschrückensteak an Rotweinjus, Rosenkohl und 20w, Butterspätzle	Gemüse-Cous-Cous mit 26 Joghurtdip
3,60 €	3,60 €	3,60 €	5,90 €	3,60 €
Fit Menü	Fit Menü	Fit Menü	Fit Menü	Fit Menü
Gratinierte Maultaschen 20(W),23,28 auf Tomatenspinat	2,4,20,25,30 Frühlingsrolle mit Süß-Saurem Dip und Basmatireis		20w,23,26Maisplinsen an gegrilltem Gemüse mit Limettensauce26	26,28 Rahmige Kartoffelsuppemit bunter Gemüseeinlage und Baguette20,23,25,26,28,29,30,3
3,30 €	3,30 €	3,30 €	3,30 €	3,30 €
1				

### OTHER, SIMILAR SOLUTIONS

- http://tabula.technology/
  - Tabula is a tool for liberating data tables locked inside PDF files.
  - Used by i.e. https://blog.openelections.net/
- Online converters
  - http://pdftables.com/ https://online2pdf.com/pdf2excel https://www.pdftoexcelonline.com/ http://www.pdf2txt.de/ etc.

# PLATOON - ALL GUTS, NO GLORY

### **LOOK INSIDE PDF**

- PDF standard is <del>old</del> mature
- Versions from 1.0 to 1.7
- Basically human-readable
- ... but rarely in the wild
- Page dimensions defined by "MediaBox"
- Unit used is "pt" (point), 1/72 inch, 0.353 mm
- Coordinate system with 0/0 at lower left corner

### MAKE IT (SLIGHTLY) READABLE

- pdftk ... uncompress
- Not really helpful with my PDF file
- QPDF does a much better job
- Still no real insight into problem domain

### **ELEMENTARY EXPERIMENTS**

- Simple 2x2 table
- With cell borders
- Save as PDF
- Look at PDF source
- Revelation: cells don't really matter

top right top left bottom left bottom right 

### **CPAN PDF MODULES**

- Way too many!
- Some are old
- Some are unmaintained
- Some are one-off forks
- Look for alternatives

# FULL METAL JACKET — THE TET OFFENSIVE

### **PDFLIB TET**

- PDFlib: PDF developer toolbox
- Decades of Postscript and PDF experience
- Bindings for Perl, lots of other languagages
- Freely available evaluation version
- As well as extensive documentation
- TET = Text and Image Extraction Toolkit

## TETML (YES, IT IS XML)

- Varying detail levels
- From word to single letters
- With bounding box etc.
- Reasonably good at detecting table structure

```
<Cell llx="109.16" llv="451.60" urx="467.89" urv="476.08">
    <Para>
        <Box llx="111.16" lly="462.28" ulx="111.16" uly="470.44" urx="604.93" ury</pre>
            <Line>
                <Text>Tomatencremesuppe Gemüsebrühe 28 Kartoffelsuppe 28</Text>
            </Line>
        </Box>
    </Para>
</Cell>
<Cell llx="554.68" llv="451.60" urx="604.93" urv="476.08">
    <Para>
        <Line>
            <Text>Zwiebelsuppe</Text>
        </Line>
    </Para>
</Cell>
<Cell llx="683.08" llv="451.60" urx="763.13" urv="476.08">
    <Para>
        <Box llx="683.08" lly="461.44" urx="763.13" ury="470.92">
            <Line>
                <Text>Erbsencremsuppe </Text>
            </Line>
            <Line>
                <Text>26,28</Text>
            </Line>
        </Box>
    </Para>
</Cell>
```

### **HOWEVER** ....

- Evaluation version
- Should not be used in production
- Features and price (Linux: 1000 €) overkill
- Absolutely recommended, just not here

# RTFM

### **CLOSER LOOK AT PDFTOTEXT (POPPLER)**

- man page could have saved me days
- crop specific area of PDF
- -- layout maintain original physical layout
- WE HAVE A WINNER!
- pdftotext --layout does exactly what we want

Donne	Mittwoch	Dienstag	⟨W 05 Montag	peiseplan KV
4. Febi	3. Februar 2016	2. Februar 2016	1. Februar 2016	
7. d oholo	Kartoffelsuppe 28	Gemüsebrühe 28	Tomatencremesuppe	Suppe
Zwiebelsu				,85 EUR
Gebackener Lebe	Gebratene Hähnchenbrust	Riesencurrywurst2,3,8	Koreanisches Feuerfleisch	
	in Rosmarin-Tomatensauce	•	vom Schwein	
mit Bohnengemüs		mit zweierlei Saucen1,2	mit Mangold	
Kartoffels	und Grillgemüse28	und Pommes Frites	und Wildreis	
	3,60 EUR	3,60 EUR	3,60 EUR	upt-
Highlight de				ichte
Hirschrückens	Hausgemachtes Schnitzel20w "Mailänder Art" mit	Kutterfrikadelle4,21	Putengulasch mit Zucchini	
Rotweinjus, Ros		Gurkensalat2,3 und	und Champignon20w,26,28 dazu	
Butterspät	Spaghetti20w und	Kartoffelpüree	Nudeln20w,23	
Fit Me	Tomatensauce 3,60 EUR Fit Menü	3,60 EUR Fit Menü	3,60 EUR Fit Menü	
Maisplinsen2 gegrilltem Ge	Gemüselasagne20w,26,28 mit	Frühlingsrolle2,4,20,25,30 mit Süß-Saurem Dip	Gratinierte Maultaschen 20(W),23,28	egeta-
	Basilikum - Tomatensauce	und Basmatireis	auf Tomatenspinat	risch

Wildreis
Beilagen Nudeln
0,70 EUR

Pommes Frites Kartoffelpüree Basmatireis Tagesgemüse28 Nussreis Spaghetti Kartoffelst Spätzle

Tagesgemüse28 Tagesgemüse28 Tagesgemüse28 Tagesgemüse2

Dessert Groß 0,80 EUR Klein 0,50 EUR Zusatzstoffe sind mit Zahlen gekennzeichnet siehe Aushang Bei Fragen zu Dussmann können Sie auf www.dussmann.com klicken!

Speiseplanänderungen vorbehalten

# THANK YOU!

- Slides will be linked on conference page
- Questions?