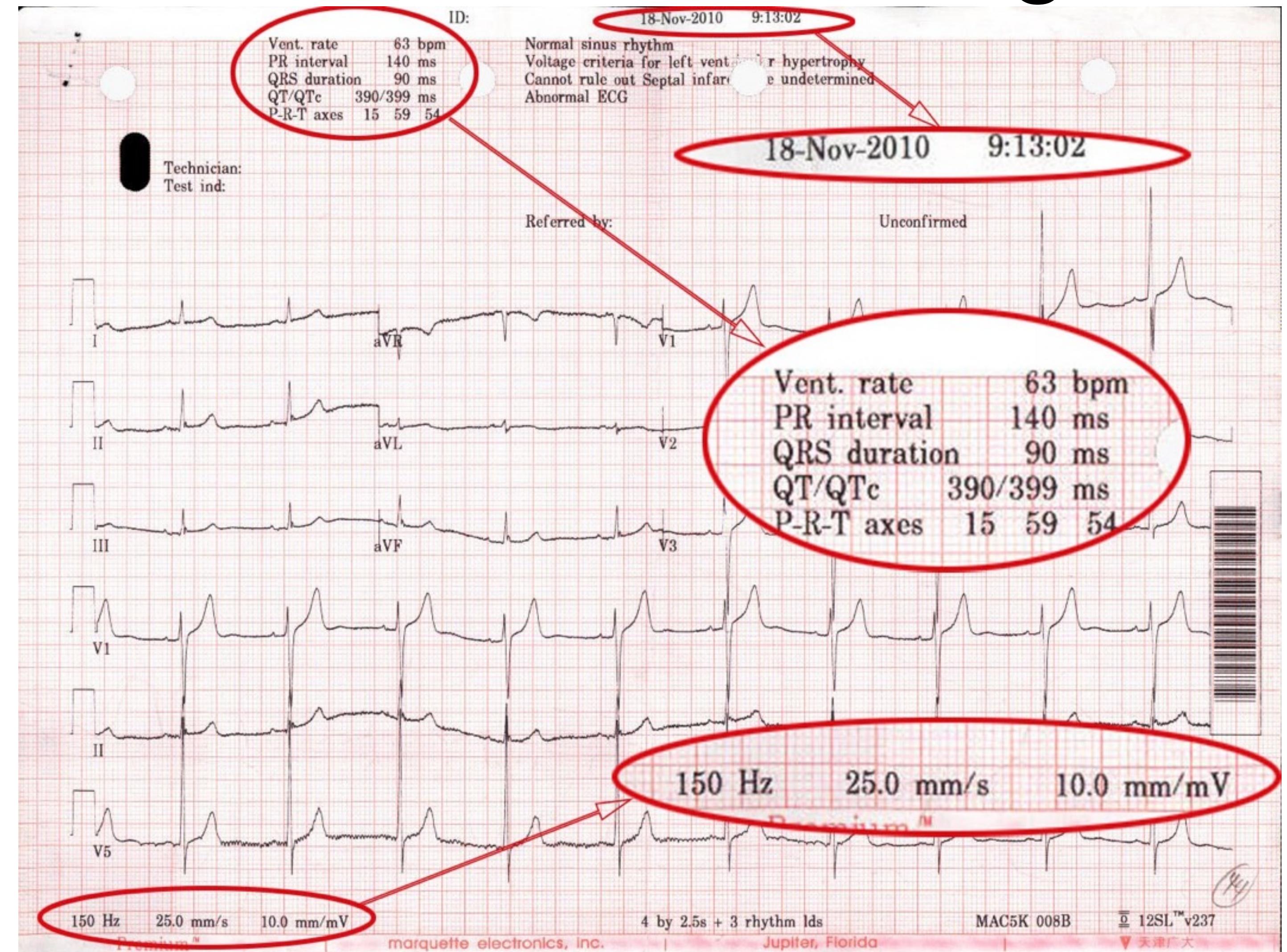


Fault-tolerant Optical Character Recognition from Semi-structured Medical Images

Jinyi Lu
JinyiLu93@gmail.com
2015/03/25

What are Semi-structured Medical Images?

- Electro-cardio-gram (ECG)



What's Fault-tolerate OCR?

- Optical Character Recognition (OCR)
 - Images of typewritten or printed text -> Machine-encoded text
- Fault-tolerate OCR

Motivation

- More than 45GB ECG images from AZ in different forms

Outline

- Input
- ODL
- Manual Correction
- Experimental Results

Outline

- Input
- ODL
- Manual Correction
- Experimental Results

Input

- XML
 - text, coordinates, layout
- Tesseract
- $\text{text_box} = \{\text{c=coor}, \text{t=v}\}$

```
<p class="ocr_par" title="box 263 33 444 119">
  <span class="ocr_l" title="box 264 33 336 45">
    <span class="ocrx_w" title="box 264 33 299 45">Vcnt.</span>
    <span class="ocrx_w" title="box 308 34 336 45">rule</span>
  </span>
  <span class='ocr_l'>
    <span class="ocrx_w" title="box 264 51 283 64">PR</span>
    <span class="ocrx_w" title="box 291 51 346 64">Interval</span>
    <span class="ocrx_w" title="box 389 52 411 64">140</span>
    <span class="ocrx_w" title="box 420 55 439 64">ms</span>
  </span>
  ...
  </span>
</p>
<p class="ocr_p" dir="ltr">
  <span class="ocr_l">
    <span class="ocrx_w" title="box 396 33 411 45">53</span>
    <span class="ocrx_w" title="box 420 33 449 48">bpm</span>
  </span>
</p>
```

Outline

- Input
- ODL
 - OCR description language (ODL)
 - Syntax
 - Semantics
- Manual Correction
- Experimental Results

Outline

- Input
- ODL
 - OCR description language (ODL)
 - Syntax
 - Semantics
- Manual Correction
- Experimental Results

Values

- All the expressions can be parsed into values

v ::=()

hskip \s

|*int*

|*float*

|*string*

|*len*

100 pixel, 100 cm

|*coor*

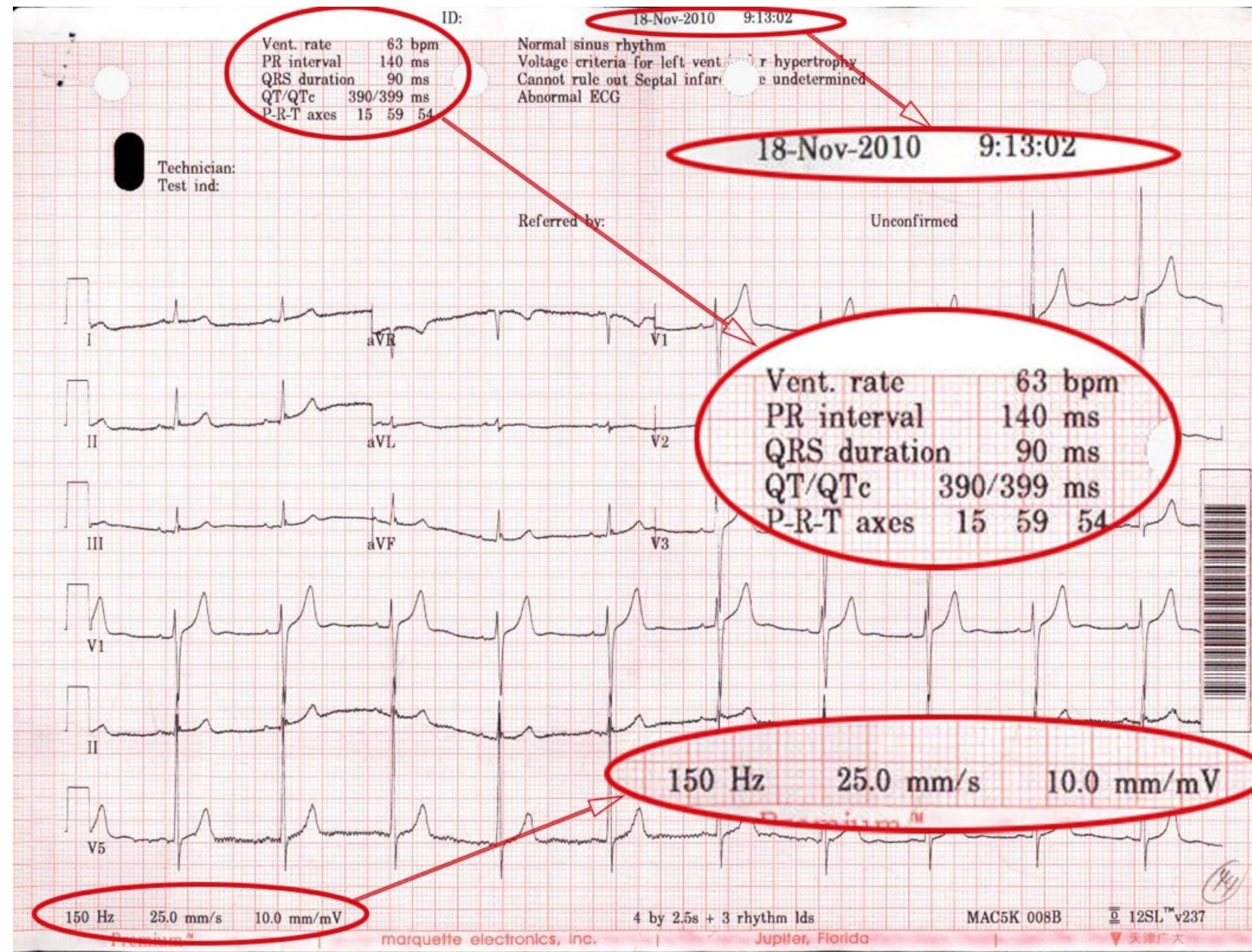
<_,0.5l, 300 pixel, 600 pixel>

|{ v_1, \dots, v_n }

Abstract Syntax

$e ::= c$	(constant)	“Vent.”
$ x$	(name)	x
$ e_0(e_1, e_2, \dots, e_n)$	(constraints)	x(int), x(<...>)
$ hskip\ e$	(horizontal skip)	hskip \s
$ vskip\ e$	(vertical skip)	vskip 10 pixel
$ \{e_1 e_2 \dots e_n\}$	(union)	{"Jan"} {"Feb"}
$ \{e_1, \dots, e_n\}$	(struct)	{"rate", x(int)}
$ e\ list$	(list)	
$ e\ as\ x$	(blinding)	
$ e_1\ bop\ e_2$	(binary operation)	

Example



Description (Abstract Syntax)

{

{

day(int , 1 , 31) ,

"-",

{

num(int , 1 , 12) |

{

"Jan" | "Feb" | "Mar" |

"Apr" | "May" | "Jun" |

"Jul" | "Aug" | "Sept" |

"Oct" | "Nov" | "Dec"

}

}

"-",

year(int)

}(<_ , _ , _ , 0.31 >) ,

...

}

Description (Surface Syntax)

```
Osource Ostruct entry_t{
    time_t(<_,_,_,0.31>)  time;
    triple_t(<0.1w,_,0.5w,_>) tri;
    inter_t(<tri.x1,tri.y0,_,_>) i;
    vskip(\n)[] skipline;
    parameter_t(<_,_,_,_>) para;
};
```

```
Ounion month_str{
    "Jan"; "Feb"; "Mar";
    "Apr"; "May"; "Jun";
    "Jul"; "Aug"; "Sept";
    "Oct"; "Nov"; "Dec";
};
```

```
Ounion month_t{
    Oint(1,12) num;
    month_str str;
};
```

```
Ostruct time_t{
    Oint(1,31) day;
    "-";
    month_t month;
    "-";
    Oint() year;
};
```

Description (Surface Syntax)

```
Ostruct triple_t{
    "Vent. rate";
    hskip(\s) skip;
    Oint(60,100) x;
    "bpm";
};

Ostruct parameter_t{
    Oint() p1;
    "Hz";
    Ofloat(3, 1) p2;
    "mm/s";
    Ofloat(3, 1) p3;
    "mm/mV";
};

Ounion inter_t{
    "Normal ECG";
    "Abnormal ECG";
};
```

Type System

$t ::=$

- $unit$
- $| int$
- $| float$
- $| len$
- $| \langle len, \ len, \ len, \ len \rangle$
- $| \{t_1 + t_2 + \dots + t_n\}$
- $| \{l_1 : t_1, \ \dots, \ l_n : t_n\}$
- $| t \ list$

$$\begin{array}{c}
 \frac{\Gamma(x) = t}{\Gamma \vdash x : t} \quad (\text{T-VARIABLE}) \\[10pt]
 \frac{\Gamma \vdash e_1 : int \ \ \Gamma \vdash e_2 : int \ \ bop \in \{+, -, *, /, \% \}}{\Gamma \vdash e_1 \ bop \ e_2 : int} \quad (\text{T-INT ARITH}) \\[10pt]
 \frac{\Gamma \vdash e_1 : int \ \ \Gamma \vdash e_2 : int \ \ bop \in \{=, !=, <, >, <=, >= \}}{\Gamma \vdash e_1 \ bop \ e_2 : bool} \quad (\text{T-INT REL}) \\[10pt]
 \frac{\Gamma \vdash e_1 : float \ \ \Gamma \vdash e_2 : float \ \ bop \in \{+, -, *, /, \% \}}{\Gamma \vdash e_1 \ bop \ e_2 : float} \quad (\text{T-FLOAT ARITH}) \\[10pt]
 \frac{\Gamma \vdash e_1 : float \ \ \Gamma \vdash e_2 : float \ \ bop \in \{=, !=, <, >, <=, >= \}}{\Gamma \vdash e_1 \ bop \ e_2 : float} \quad (\text{T-FLOAT REL}) \\[10pt]
 \frac{\Gamma \vdash e_0 : t_0}{\Gamma \vdash e_0(e_1, e_2, \dots, e_n) : t_0} \quad (\text{T-CONSTRAINT}) \\[10pt]
 \frac{\Gamma \vdash e : len}{\Gamma \vdash hskip \ e : unit} \quad (\text{T-HSKIP}) \\[10pt]
 \frac{\Gamma \vdash e : len}{\Gamma \vdash vskip \ e : unit} \quad (\text{T-VSKIP}) \\[10pt]
 \frac{for \ each \ i \ \ \Gamma \vdash e_i : t_i}{\Gamma \vdash \{e_1 | \dots | e_n\} : t_1 + \dots + t_n} \quad (\text{T-UNION}) \\[10pt]
 \frac{for \ each \ i \ \ \Gamma \vdash e_i : t_i}{\Gamma \vdash \{e_1, \dots, e_n\} : t_1 * \dots * t_n} \quad (\text{T-STRUCT}) \\[10pt]
 \frac{\Gamma \vdash e : t}{\Gamma \vdash e \ list : t \ list} \quad (\text{T-LIST})
 \end{array}$$

Outline

- Input
- ODL
 - OCR description language (ODL)
 - Syntax
 - Semantics
- Manual Correction
- Experimental Results

Semantics

- Automatically generate the most reliable parsing results
 - Input data and parsing results
 - Scoring policy
 - Evaluation rules

Input Data and Parsing Results

`text_box ::= {c = coor, t = v}`

| {}

`parse_tree ::= ((e, text_box), [parse_tree1, ..., parse_treen])`

| ()

`input_data ::= {text_box | \forall text_box in XML}`

time			
day	F	"18"	
"_"	F		
month			
str	F	"Nov"	
"_"	F		
year	F	"2010"	
triple			
"Vent. rate"	E	"Vcnt. rule"	
skip1	F		
x	E	"53"	
"bpm"	F		
skip2	F		
inter			
"Abnormal ECG"	F		
para			
p1	E	"150"	
"Hz"	F		
p2	F	"25.0"	
"mm/s"	F		
p3	E	"10.0"	
"mm/mV"	F		

Scoring Policy

- Score for how likely a text box is matched with the expression
- Why?
 - Errors and noises in the OCR results
 - Rough description
- Two kinds
 - Data description
 - Spatial description

Data description

- To tolerate the errors and noises in the OCR results

$$es(c, t) = ed(c, t)$$

$$es(x(int, a, b), t) = \min\{ed(i, t) \mid \forall i \in Z, i \in [a, b]\}$$

$$es(x(float, l, p, a, b), t) = \min\{ed(i, t) \mid \forall i \in R', i \in [a, b]\}$$

Data description

“Vent. rate”	“Vcnt. rule”	3
x(int, 60, 100)	53	1
y(float, 3, 1)	10.0	2

Scoring Policy

- Spatial description
 - Spatial description is used to describe the rough spatial relationship based on human estimates
 - cc: the position of the current cursor

$$ss(cc, coor) = ||cc - coor||$$

Scoring Policy

$$\text{score}(e, tb, cc) = \text{es}(e, tb.t) + k * \text{ss}(cc, tb.c)$$

Judgment Form

- E: environment
 - cc: the position of the current cursor
- D: input data
- D': remaining data

Judgment Form : $E, D; e \Downarrow (D'; \text{parse_tree})\text{list}$

Judgment Form

- Why list?
 - Scoring policy: lots of satisfied text box
 - All candidate parsing results
 - Final result: sum up the scores of all leaf nodes

Evaluation Rules

- Examples
 - {"Vent. rate", x(int, 60, 100)}
 - {{<...>, "Vcnt. rule"}, {<...>, "b"}, {<...>, "53"}}
- Try all the candidates

$$\overline{E, \{\}; x \Downarrow []}$$

(E-EMPTY)

$$\frac{p \in D \ E(c) = coor \ Cons(p, coor, const) = true \ E, (D - p); x \Downarrow r}{E, D; const \Downarrow ((D - p); (const, p), []) :: r} \quad (\text{E-C1})$$

$$\frac{p \in D \ E(c) = coor \ Cons(p, coor, const) = false \ E, (D - p); x \Downarrow r}{E, D; const \Downarrow (D; (const, \{}), \[])) :: r} \quad (\text{E-C2})$$

$$\frac{p \in D \ E(c) = coor \ E, (D - p); x \Downarrow r}{E, D; x \Downarrow ((D - p); (x, p), \[])) :: r} \quad (\text{E-X})$$

$$\frac{E, D; e_0 \Downarrow r_0 \quad \forall i : r_{i+1} = Filter(r_i, e_i)}{E, D; e_0(e_1, e_2, \dots, e_n) \Downarrow r_{n+1}} \quad (\text{E-CONSTRAINT})$$

$$\frac{Find(D, coor, nil) = D' \ E[c \mapsto Begin(D')], D'; e \Downarrow (D''; t)list}{E, D; e(coor) \Downarrow (D - D' + D''; t)list} \quad (\text{E-COOR})$$

$$\frac{E, D; e \Downarrow r \ E(c) = coor \ Hskip(D, coor, r) = D'}{E, D; hskip \ e \Downarrow (D'; \[])) :: \[]} \quad (\text{E-HSKIP})$$

$$\frac{E, D; e \Downarrow r \ E(c) = coor \ Vskip(D, coor, r) = D'}{E, D; vskip \ e \Downarrow (D'; \[])) :: \[]} \quad (\text{E-VSKIP})$$

$$\frac{E, D; e_1 \Downarrow r_1 \ E, D; \{e_2|..|e_n\} \Downarrow r_2}{E, D; \{e_1|e_2|..|e_n\} \Downarrow r_1 + r_2} \quad (\text{E-UNION})$$

$$\frac{E, D; e_1 \Downarrow (E', D'; parse_tree)list \quad \forall i : E'_i[c \mapsto End(D - D'_i)], D'_i; \{e_2, \dots, e_n\} \Downarrow r_i}{E, D; \{e_1, e_2, \dots, e_n\} \Downarrow \sum_{i=1}^m (D'; parse_tree)list * r_i} \quad (\text{E-STURCT})$$

$$\frac{E, D; head \ e \ list \Downarrow (E', D'; parse_tree)list \quad \forall i : E'_i[c \mapsto End(D - D'_i)], D'_i; tail \ e \ list \Downarrow r_i}{E, D; e \ list \Downarrow \sum_{i=1}^m (D'; parse_tree)list * r_i} \quad (\text{E-LIST})$$

Outline

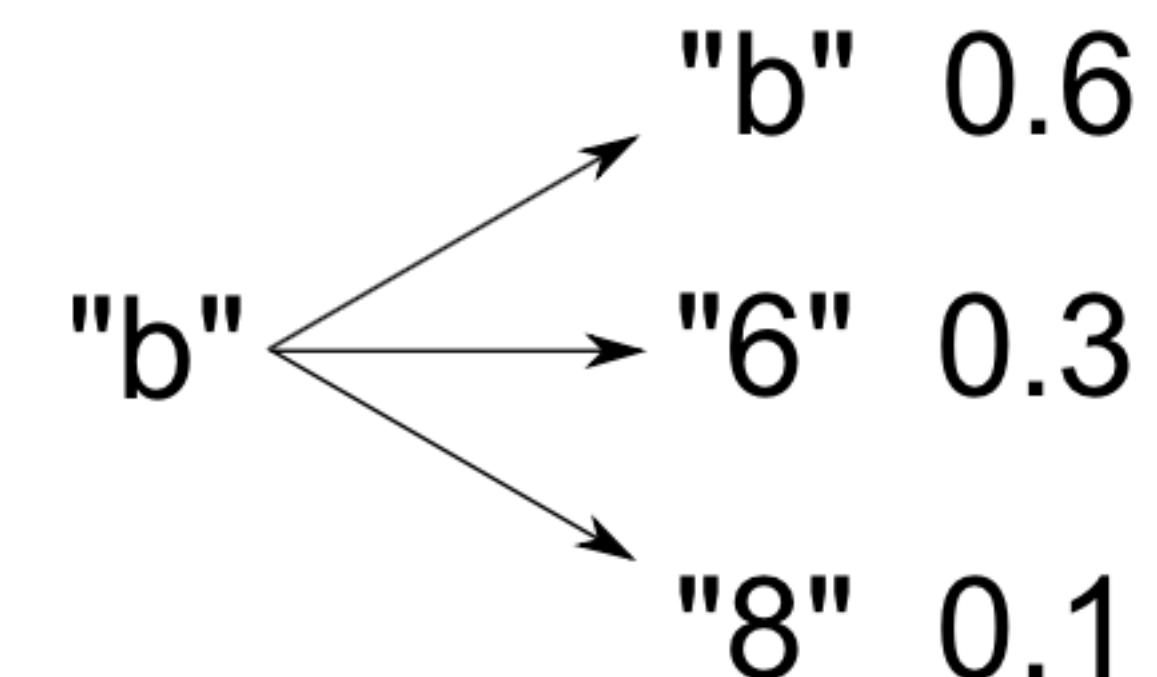
- Input
- ODL
- Manual Correction
 - Incremental Learning Correction Model
 - Manual Correction Policy
- Experimental Results

Correction Model

- Correction strategies S
 - m: original substring
 - n: candidate substring after correction
 - A: set of the correct candidates and probabilities

$$S = \{(m, n, p) | \forall (n, p) \in A\}$$

- Correction model M



New Score Policy

$$cor(t, S) = \{(t', p) | \forall (m, n, p) \in S, rep(t, m, n) = t'\}$$

$$esm(e, t, M) = \min\{es(e, t') * p | \forall S \in M, \forall (t', p) \in cor(t, S)\}$$

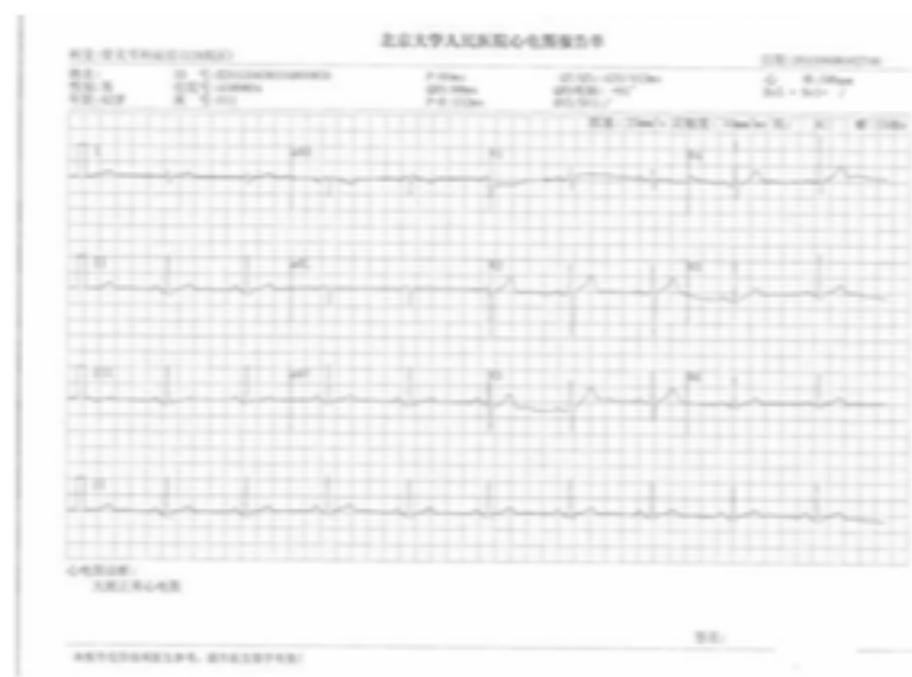
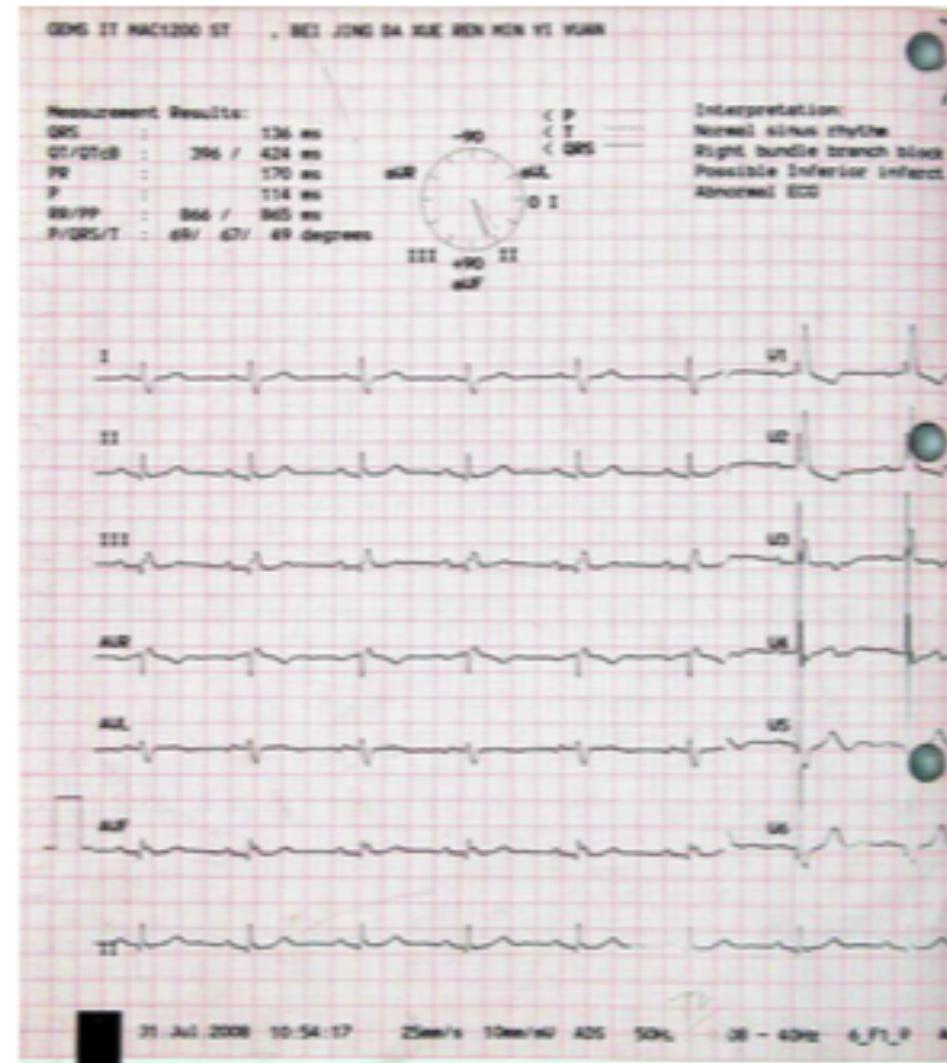
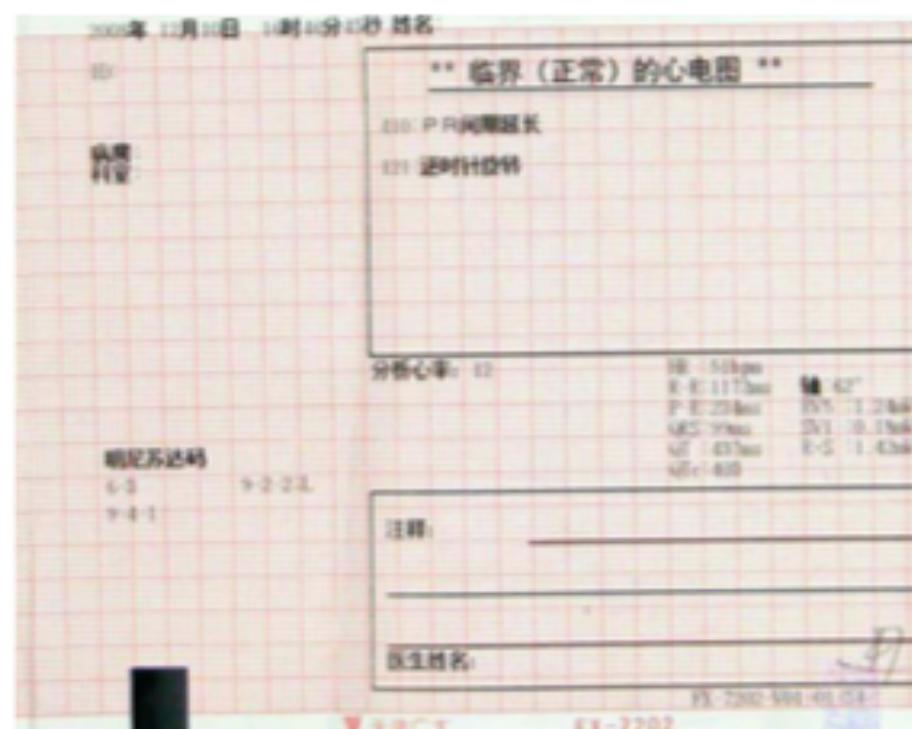
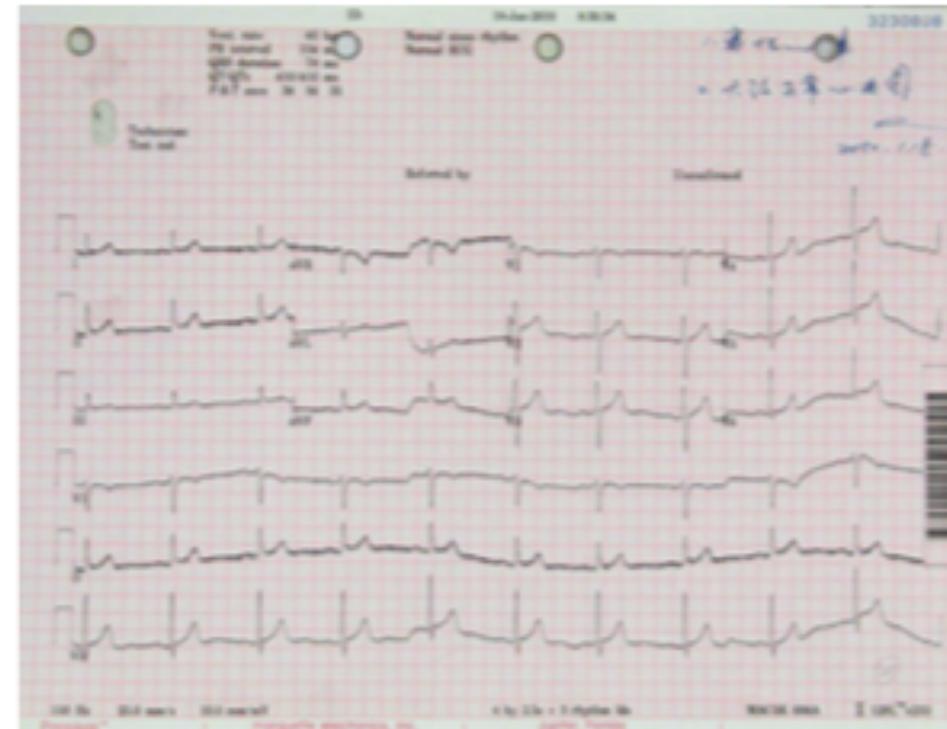
Manual Correction Policy

- Random
- Most frequent error description elements

Outline

- Input
- ODL
- Manual Correction
- Experimental Results
 - Dataset
 - Other Approach
 - Extraction Accuracy
 - Incremental Manual Correction

Dataset



Format	1	2	3	4
Number of Images	124	113	102	97
Number of Attributes per Image	17	16	18	15

Other Approach

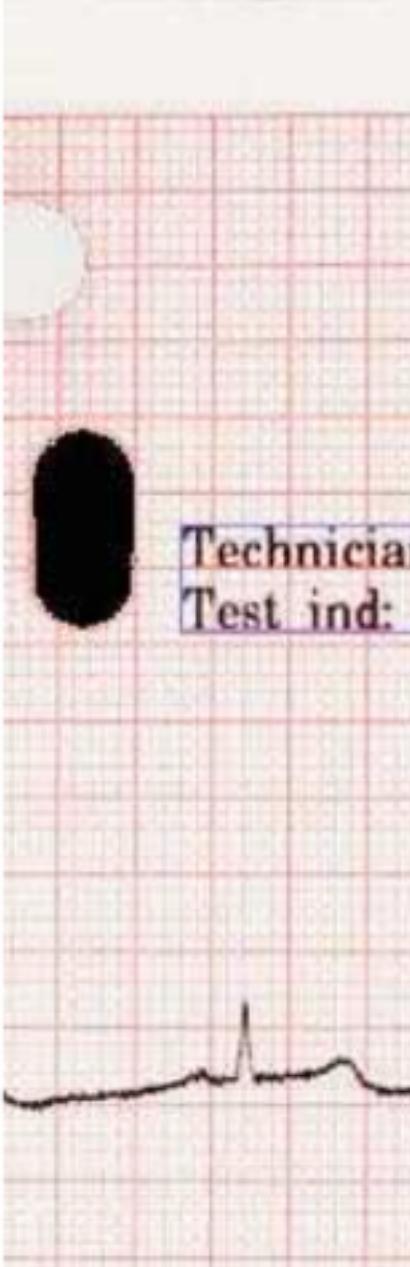
- Exact match
- Zonal OCR
- Page layout analysis

ID:	18-Nov
Vent. rate	63 bpm
PR interval	140 ms
QRS duration	90 ms
QT/QTc	390/399 ms
P-R-T axes	15 59 54
Technician:	
Test ind:	
Referred by:	

Normal sinus rhythm
Voltage criteria for left
Cannot rule out Septal
Abnormal ECG

Other Approach

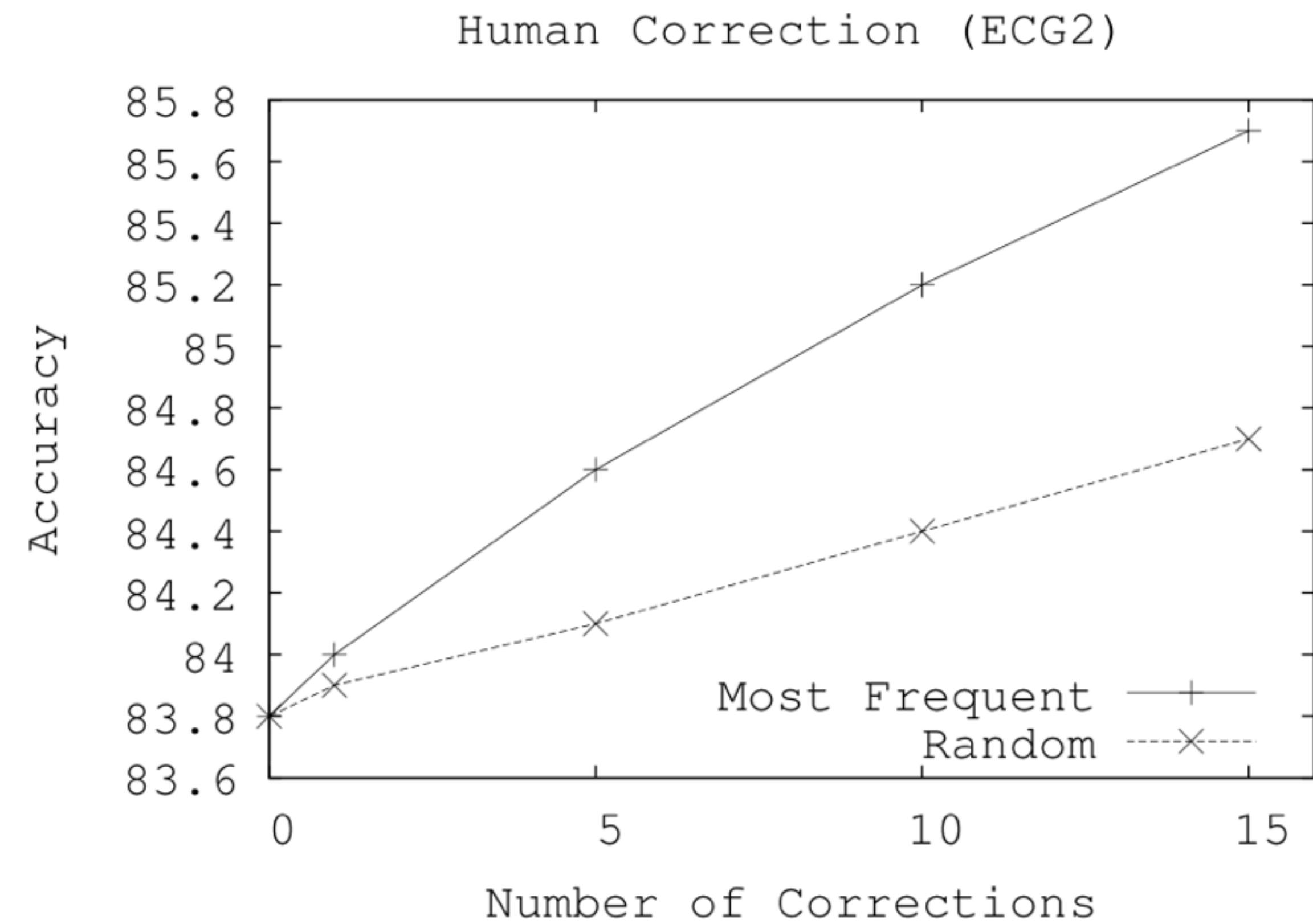
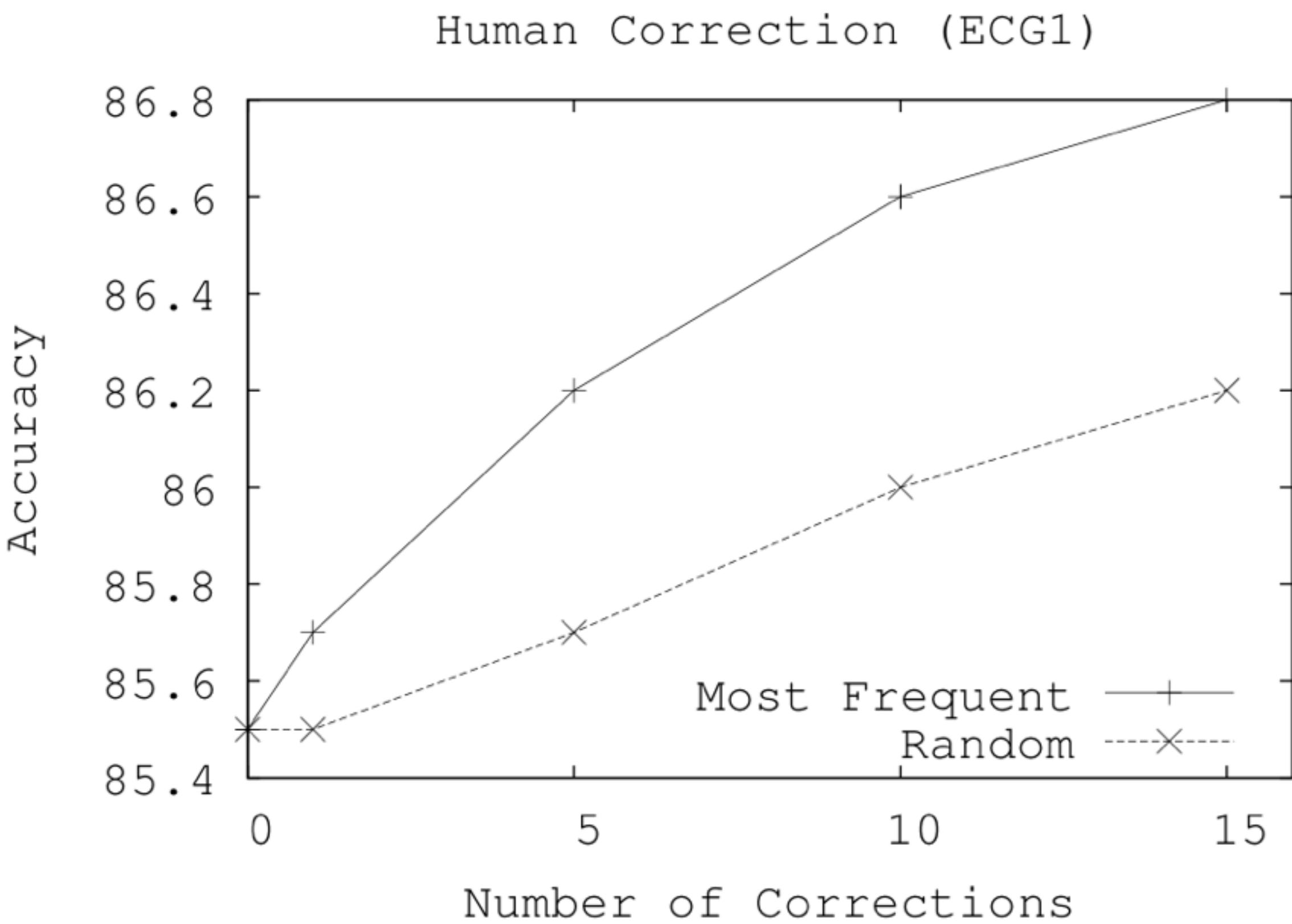
- Exact match
- Zonal OCR
- Page layout analysis

		ID: 18-Nov-2010 9:13:02										
<table border="1"><tr><td>Vent. rate</td><td>63 bpm</td></tr><tr><td>PR interval</td><td>140 ms</td></tr><tr><td>QRS duration</td><td>90 ms</td></tr><tr><td>QT/QTc</td><td>390/399 ms</td></tr><tr><td>P-R-T axes</td><td>15 59 54</td></tr></table>		Vent. rate	63 bpm	PR interval	140 ms	QRS duration	90 ms	QT/QTc	390/399 ms	P-R-T axes	15 59 54	Normal sinus rhythm Voltage criteria for left ventricular hypertrophy Cannot rule out Septal infarction Abnormal ECG
Vent. rate	63 bpm											
PR interval	140 ms											
QRS duration	90 ms											
QT/QTc	390/399 ms											
P-R-T axes	15 59 54											
Technician: Test ind:		Referred by:										
												
<table border="1"><tr><td>Vent. rate</td><td>74 bpm</td></tr><tr><td>PR interval</td><td>134 ms</td></tr><tr><td>QRS duration</td><td>78 ms</td></tr><tr><td>QT/QTc</td><td>380/421 ms</td></tr><tr><td>P-R-T axes</td><td>64 58 39</td></tr></table>		Vent. rate	74 bpm	PR interval	134 ms	QRS duration	78 ms	QT/QTc	380/421 ms	P-R-T axes	64 58 39	Normal sinus rhythm Normal ECG
Vent. rate	74 bpm											
PR interval	134 ms											
QRS duration	78 ms											
QT/QTc	380/421 ms											
P-R-T axes	64 58 39											
Technician: Test ind:		Referred by:										

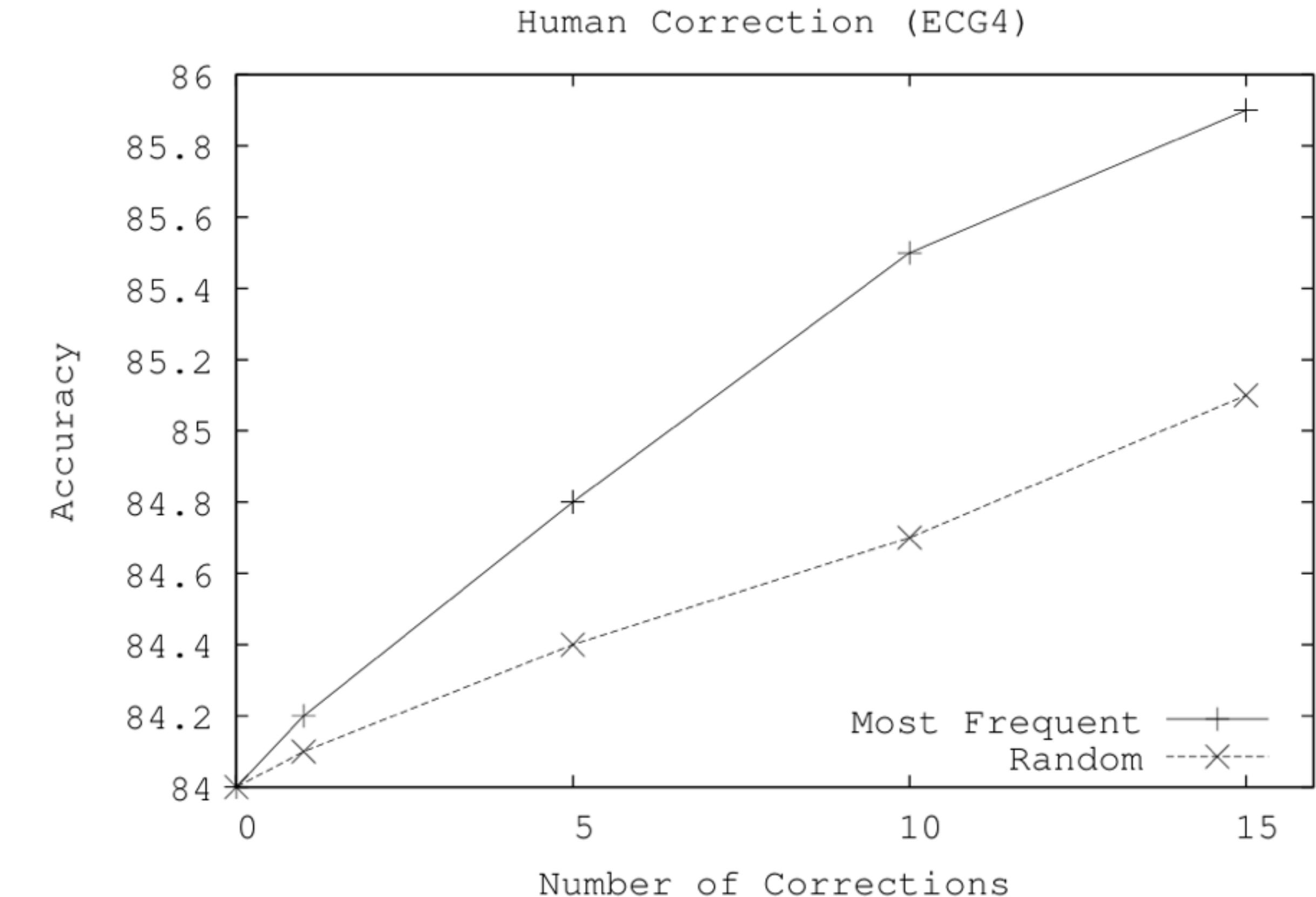
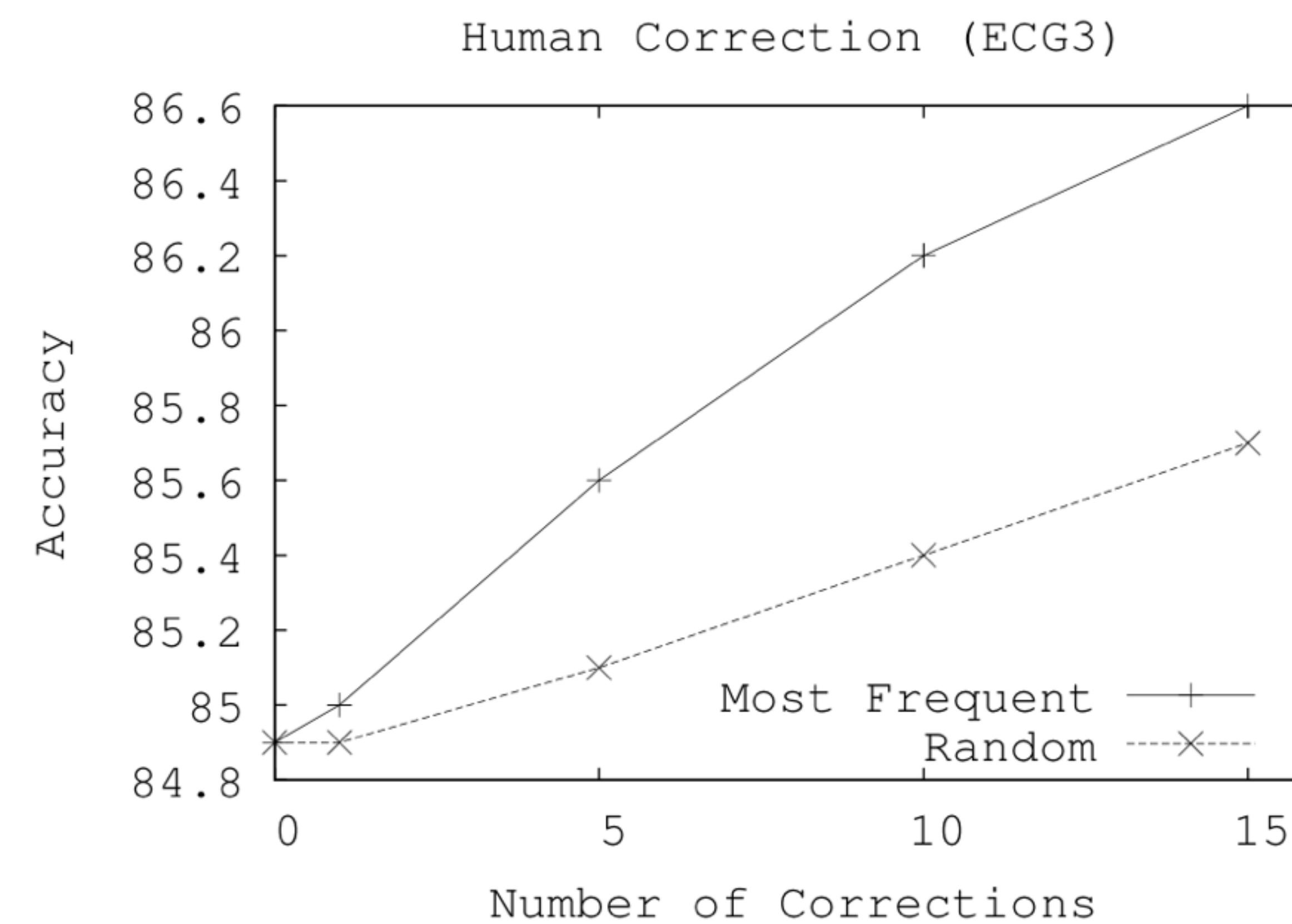
Extraction Accuracy

Format	1	2	3	4
Exact Match	58.8%	56.3%	61.1%	53.4%
Zonal OCR	81.2%	79.8%	81.7%	80.6%
Page Layout	79.7%	80.2%	81.2%	81.1%
Our Fuzzy Match	85.5%	83.8%	84.9%	84.0%

Incremental Manual Correction



Incremental Manual Correction



Thank you!