

# PERFORMANCE MEASURES AND EVALUATION

- Estimating performance
- Classic performance metrics
- Performance metrics
- How well are ST episodes detected?
- Characteristics of transient ST episodes
- Estimating performance of ST episode detection
- Estimating performance of detecting total ST episode time
- Estimating performance of measuring the ST segment deviation
- Estimating performance of classifying ischaemic and non-ischaemic heart - rate related ST episodes
- Assessing the robustness (predicting the “real world” performance)

# Estimating performance

- **Development and testing database ??**
- In the field of developing analyzers of biomedical signals **use development database and report performance achieved**, and possibly **try to predict the analyzer's performance in the “real world”** (ANSI/AAMI, Us. FDA, Physionet)  
(Examples: MIT/BIH DB, QT DB, ESC DB, LTST DB, TPEHG DB)
- **Due to desired comparison of the results of the evaluation**, it is necessary to use **ALL records** of the database available and **UNIQUE performance metrics**

# Classic performance measures

- Performance evaluation

## Classic event oriented performance matrix

		Analyzer	Analyzer
		EVENT	NON-EVENT
Reference	event	<i>TP</i>	<i>FN</i>
Reference	non-event	<i>FP</i>	<i>(TN)</i>

*TP* – number of correctly detected events

*FN* – number of missed events

*FP* – number of falsely detected events

*TN* – number of correctly rejected non-events  
(undefined for detection task !)

# Classic performance measures

- Performance evaluation (detection task)

		Analyzer	Analyzer
		EVENT	NON-EVENT
Reference	event	<i>TP</i>	<i>FN</i>
Reference	non-event	<i>FP</i>	<i>(TN)</i>

Sensitivity:

$$Se = \frac{TP}{TP + FN}$$

The proportion of EVENTS which were correctly detected as events

Positive predictivity:

$$+P = \frac{TP}{TP + FP}$$

The proportion of detected EVENTS which actually were events

# Classic performance measures

- Performance evaluation (classification task)

		Analyzer	Analyzer
		EVENT	NON-EVENT
Reference	event	<i>TP</i>	<i>FN</i>
Reference	non-event	<i>FP</i>	<i>TN</i>

**Sensitivity:**

$$Se = \frac{TP}{TP + FN}$$

The proportion of EVENTS which were correctly classified as events

**Positive predictivity:)**

$$+P = \frac{TP}{TP + FP}$$

The proportion of classified EVENTS which actually were events

**Specificity:**

$$Sp = \frac{TN}{TN + FP}$$

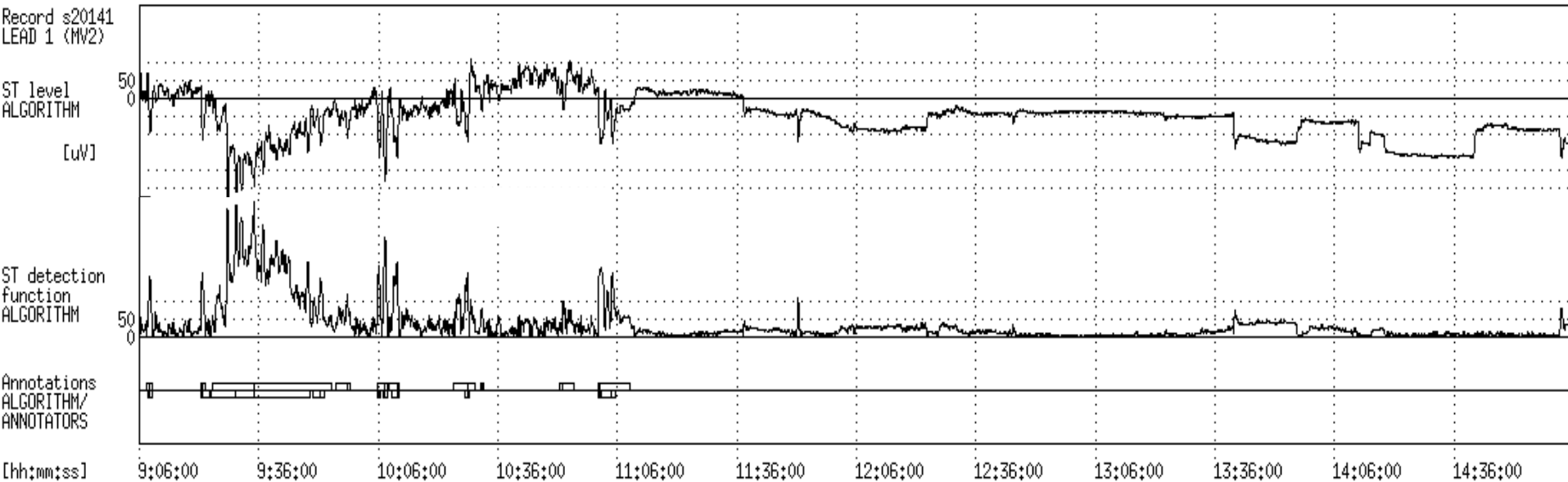
The proportion of NON-EVENTS which were correctly classified as non-events)

# Performance metrics

- The evaluation of an ST detection algorithm or analyzer should answer the following questions:
  - **How well are ST episodes detected ?**
  - How reliably is ST episode or ischaemic ST episode duration measured ?
  - How accurately are ST deviations measured ?
  - How well are ischaemic and non-ichaemic heart rate related ST episodes differentiated ?
  - How well will the ST analyzer perform in the “real world” ?

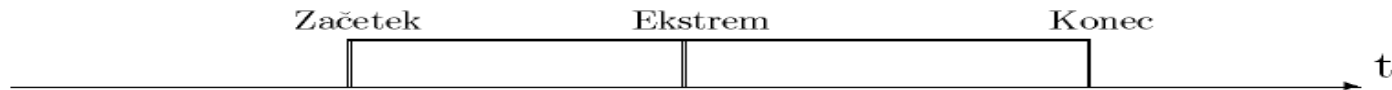


# How well are ST episodes detected?

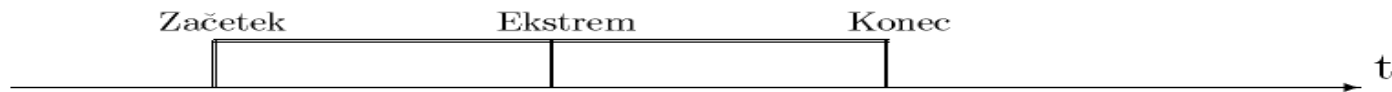


# Characteristics of transient ST episodes

Analizatorjeva epizoda ST



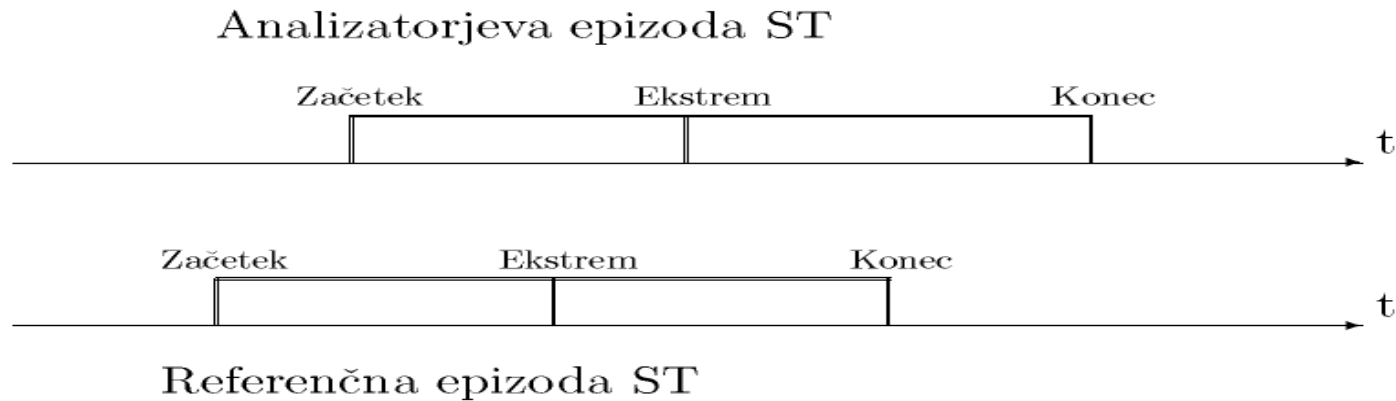
Referenčna epizoda ST



- Number
- Type (ischaemic, due to heart frequency changes)
- Length
- ST segment level at extrema deviation



# Estimating performance of ST episode detection



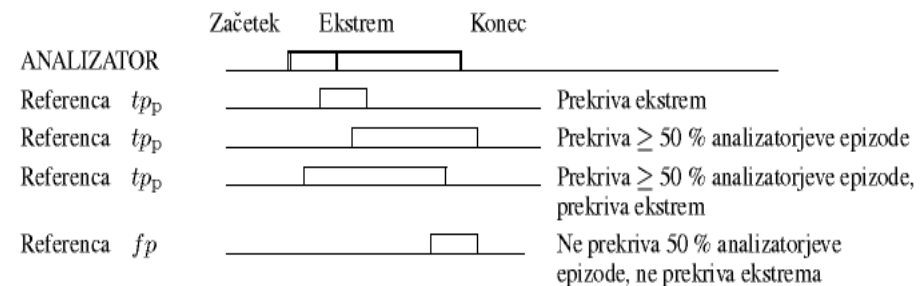
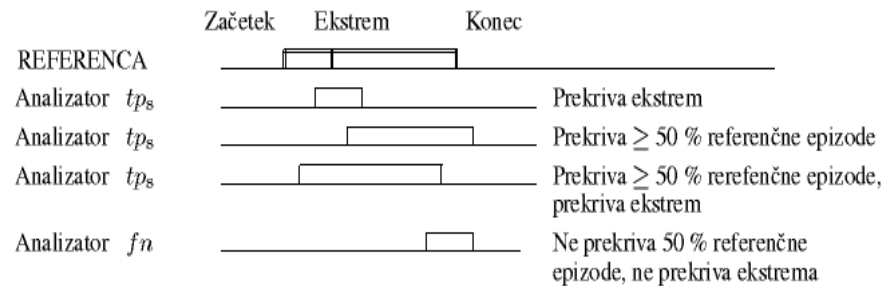
## Characteristics of analyzer-annotated and reference episodes:

- Episodes contain time dimension
- They differ in length
- There is no one-to-one correspondence
- There is no “non-events”

# Estimating performance of ST episode detection

## Estimating performance of transient ST segment episode detection

TPs, TPp, FN, FP: **Matching test** for reference and analyzer-annotated episodes



TPs - Reference episodes that satisfied the matching criteria

TPp - Analyzer-annotated episodes that satisfied the matching criteria

FN - Missed episodes

FP - Falsely detected episodes

# Estimating performance of ST episode detection

Estimating performance of transient ST segment episode detection

Performance measures (ANSI/AAMI, Us. FDA)

$Se$		Analyzer	
		ST epi.	Other
Ref.	ST epi.	$TPs$	$FN$
	Other	-	-

## $Se$ - Sensitivity

An estimate of the likelihood of detecting an ST episode

$$Se = TPs / (TPs + FN)$$

$+P$		Analyzer	
		ST epi.	Other
Ref.	ST epi.	$TPp$	-
	Other	$FP$	-

## $+P$ – Positive predictivity

An estimate of the likelihood that a detection is a true ST episode:

$$+P = TPp / (TPp + FP)$$

# Estimating performance of ST episode detection

(LTST DB, Protocol B, 908 combined ST segment episodes)  
(Time domain and KLT approach)

Se		Algorithm	
		ST epis	Other
Ref	ST epis	723	185
Ref	Other	-	-

+P		Algorithm	
		ST epis	Other
Ref	ST epis	750	-
Ref	Other	208	-

ST epis – Combined ST segment episodes

Se = 79.6%      +P = 78.3%

# Estimating performance of ST episode detection

<i>Technique</i>		<i>ESC DB</i>	
		<i>SE [%]</i>	
		<i>Se</i>	<i>+P</i>
Time domain	[g]	81	76
	[a]	84	81
RMS method	[g]	–	–
	[a]	84.7	86.1
Time domain	[g]	79.2	81.4
	[a]	81.5	82.5
KLT approach	[g]	85.2	86.2
	[a]	87.1	87.7
Time domain, KLT	[g]	77.2	86.3
	[a]	81.3	89.2
Neural net	[g]	85.0	68.7
	[a]	88.6	78.4
Neural net, KLT	[g]	–	–
	[a]	77	86



# Estimating performance of ST episode detection

<i>Technique</i>		<i>ESC DB</i>		<i>LTST DB</i>		
		<i>SE [%]</i>		<i>SE [%]</i>		
		<i>Se</i>	<i>+P</i>		<i>Se</i>	<i>+P</i>
Time domain	[g]	81	76		-	-
	[a]	84	81		-	-
RMS method	[g]	-	-		-	-
	[a]	84.7	86.1		-	-
Time domain	[g]	79.2	81.4		-	-
	[a]	81.5	82.5		-	-
KLT approach	[g]	85.2	86.2	->	77.0	58.8
	[a]	87.1	87.7	->	74.0	61.4
Time domain, KLT	[g]	77.2	86.3	<-	79.6	78.3
	[a]	81.3	89.2	<-	78.9	80.7
Neural net	[g]	85.0	68.7		-	-
	[a]	88.6	78.4		-	-
Neural net, KLT	[g]	-	-		-	-
	[a]	77	86		-	-



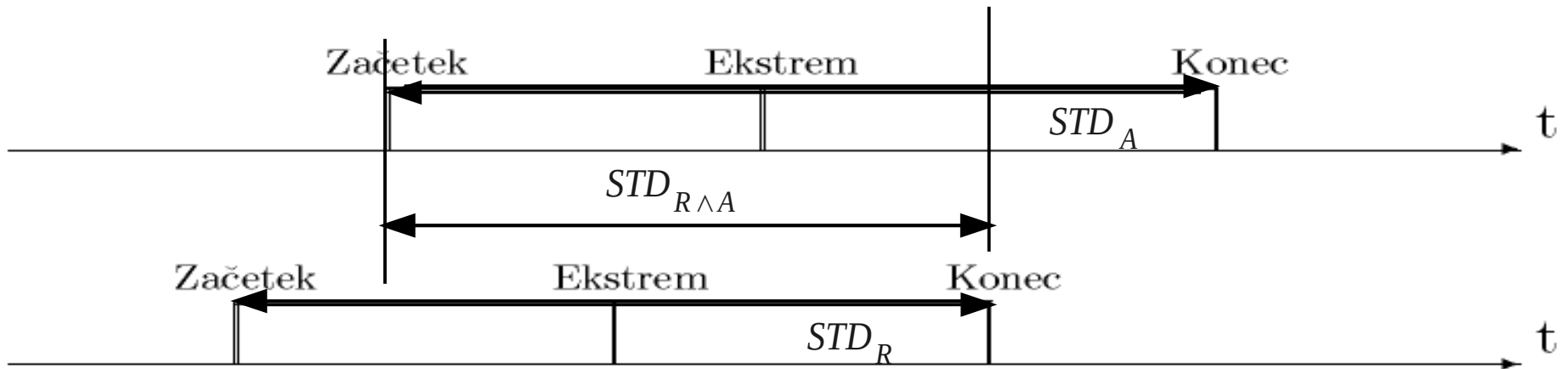


# Performance metrics

- The evaluation of an ST detection algorithm or analyzer should answer the following questions:
  - How well are ST episodes detected ?
  - **How reliably is ST episode or ischaemic ST episode duration measured ?**
  - How accurately are ST deviations measured ?
  - How well are ischaemic and non-ichaemic heart rate related ST episodes differentiated ?
  - How well will the ST analyzer perform in the “real world” ?

# Estimating performance of detecting total ST episode time

## Analizatorjeva epizoda ST



## Referenčna epizoda ST

$$STD \ Se = \frac{STD_{R \wedge A}}{STD_R} \quad STD + P = \frac{STD_{R \wedge A}}{STD_A}$$

(ANSI/AAMI, Us. FDA)



# Estimating performance of detecting total ST episode time

<i>Technique</i>		<i>ESC DB</i>				<i>LTST DB (Protocol B)</i>				
		<i>SE [%]</i>		<i>SD [%]</i>		<i>SE [%]</i>		<i>SD [%]</i>		
		<i>Se</i>	<i>+P</i>	<i>Se</i>	<i>+P</i>	<i>Se</i>	<i>+P</i>	<i>Se</i>	<i>+P</i>	
Time domain	[g]	81	76	-	-	-	-	-	-	
	[a]	84	81	-	-	-	-	-	-	
RMS method	[g]	-	-	-	-	-	-	-	-	
	[a]	84.7	86.1	75.3	68.2	-	-	-	-	
Time domain	[g]	79.2	81.4	-	-	-	-	-	-	
	[a]	81.5	82.5	-	-	-	-	-	-	
KLT approach	[g]	85.2	86.2	75.8	78.0	->	77.0	58.8	48.5	47.8
	[a]	87.1	87.7	78.2	74.1	->	74.0	61.4	54.8	58.4
→ Time domain, KLT	[g]	77.2	86.3	67.5	69.2	<-	79.6	78.3	68.4	67.3
	[a]	81.3	89.2	77.6	68.9	<-	78.9	80.7	73.1	74.9
Neural net	[g]	85.0	68.7	73.0	69.5	-	-	-	-	
	[a]	88.6	78.4	72.2	67.5	-	-	-	-	
Neural net, KLT	[g]	-	-	-	-	-	-	-	-	
	[a]	77	86	-	-	-	-	-	-	

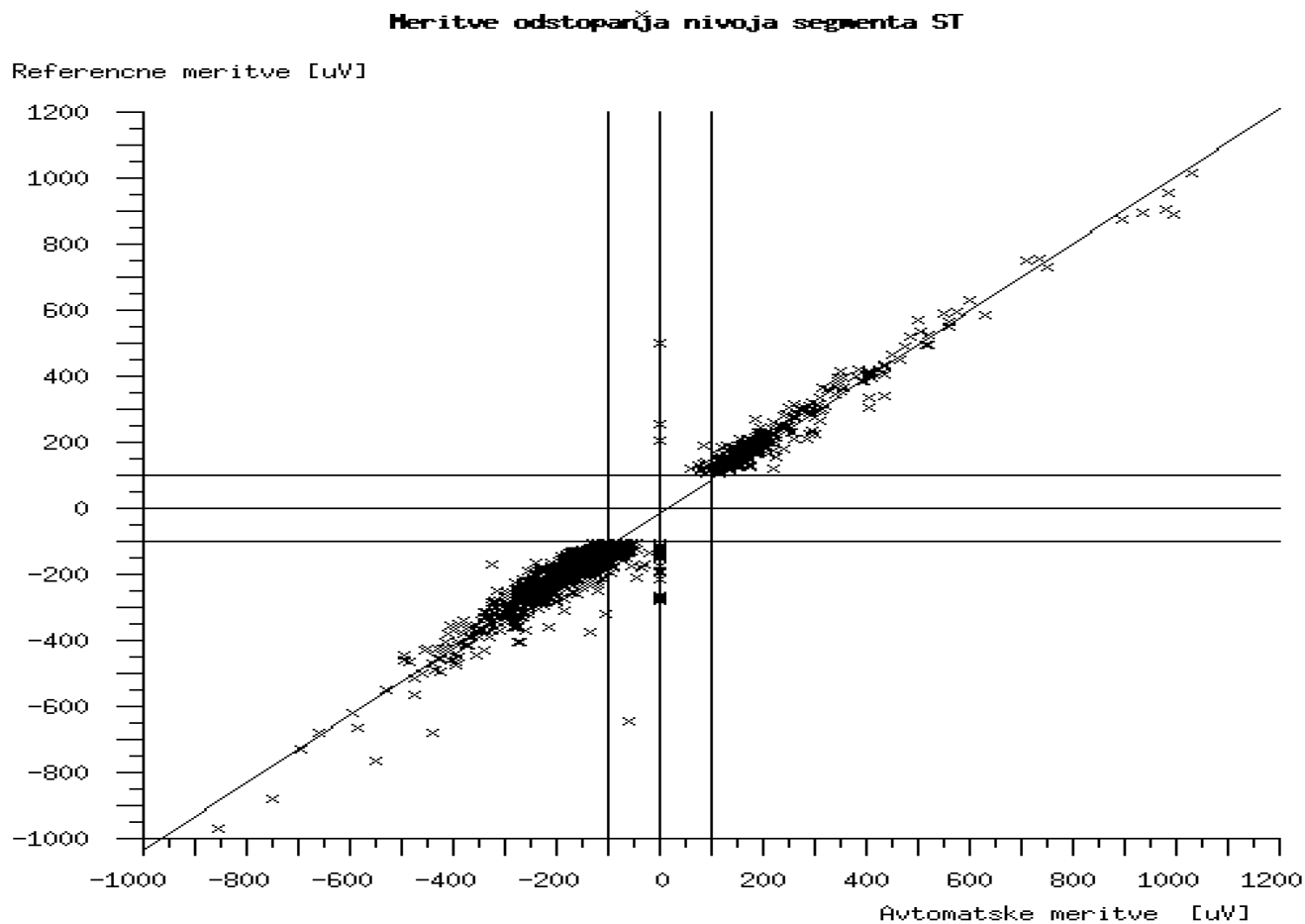
# Performance metrics

- The evaluation of an ST detection algorithm or analyzer should answer the following questions:
  - How well are ST episodes detected ?
  - How reliably is ST episode or ischaemic ST episode duration measured ?
  - **How accurately are ST deviations measured ?**
  - How well are ischaemic and non-ichaemic heart rate related ST episodes differentiated ?
  - How well will the ST analyzer perform in the “real world” ?

# Estimating performance of measuring the ST segment deviation

- **Measurement error = Measurement of the analyzer - Reference measurement**
- Mean error
- Standard deviation
- Correlation coefficient
- Regression line
- Percentage of measurements of which error exceeded  $100\mu\text{V}$   
(Discrepant ST segment measurement percentage,  $p(100\mu\text{V})$ )
- Measurement error in  $\mu\text{V}$  which was not exceeded by 95% of measurements  
(Discrepant ST segment deviation measurement value,  $e(95\%)$ )

# Estimating performance of measuring the ST segment deviation



<b>Število epizod</b> =	<b>1364</b>	<b>Povprečje [uV]</b> =	<b>16.41</b>
<b>St. dev. [uV]</b> =	<b>63.13</b>	<b>Korel. koef.</b> =	<b>0.96</b>
<b>p(100uV) [%]</b> =	<b>4.25</b>	<b>e(95%) [uV]</b> =	<b>90.00</b>
<b>Ref. [uV]</b> =	<b>1.02</b>	<b>Meritev - 14.36 [uV]</b>	

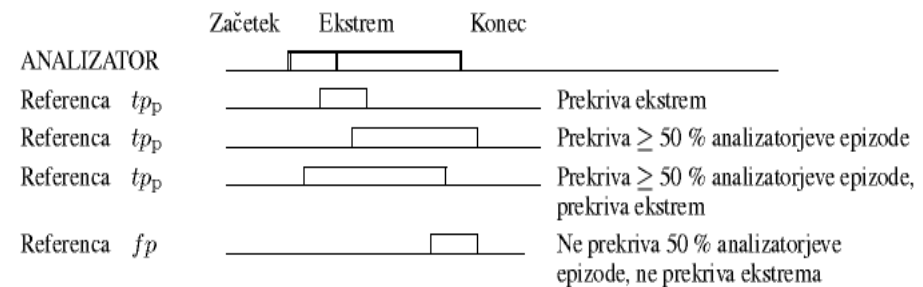
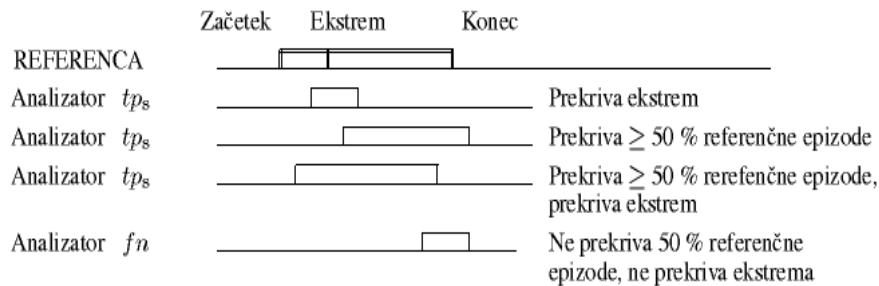
# Performance metrics

- The evaluation of an ST detection algorithm or analyzer should answer the following questions:
  - How well are ST episodes detected ?
  - How reliably is ST episode or ischaemic ST episode duration measured ?
  - How accurately are ST deviations measured ?
  - **How well are ischaemic and non-ichaemic heart rate related ST episodes differentiated ?**
  - How well will the ST analyzer perform in the “real world” ?

# Estimating performance of classifying ischaemic and non-ischaemic heart-rate related ST episodes

## Estimating performance of transient ST segment episode detection

TPs, TPp, FN, FP: **Matching test** for reference and analyzer-annotated episodes



TPs - Reference episodes that satisfied the matching criteria

TPp - Analyzer-annotated episodes that satisfied the matching criteria

FN - Missed episodes

FP - Falsely detected episodes

# Estimating performance of classifying ischaemic and non-ischaemic heart-rate related ST episodes

- **Extended matching test** ( $SR(u)$ ,  $SA(v)$  – status of reference or analyzer-annotated episode after the matching test)

*if match with ischemic analyzer-annotated ST episodes then*

*$S_R(u) = ischemic;$*

*else if match with heart rate related analyzer-annotated ST episodes then*

*$S_R(u) = heart\ rate\ related;$*

*else*

*$S_R(u) = missed;$*

*endif*

*if match with ischemic reference-annotated ST episodes then*

*$S_A(v) = ischemic;$*

*else if match with heart rate related reference-annotated ST episodes then*

*$S_A(v) = heart\ rate\ related;$*

*else*

*$S_A(v) = falsely\ detected;$*

*endif*

# Estimating performance of classifying ischaemic and non-ischaemic heart-rate related ST episodes

Se		Analyzer		
		Isch	Hrr	Other
Ref.	Isch	a	b	c
	Hrr	d	e	f
	Other	-	-	-

+P		Analyzer		
		Isch	Hrr	Other
Ref.	Isch	g	h	-
	Hrr	i	j	-
	Other	k	l	-

Isch - Ischaemic ST segment episodes

Hrr - Heart rate related ST segment episodes

Other - No ST segment changes

ST epi. - ST segment episodes (Isch + Hrr)



# Estimating performance of classifying ischaemic and non-ischaemic heart-rate related ST episodes

Se		Analyzer		
		Isch	Hrr	Other
Ref.	Isch	a	b	c
	Hrr	d	e	f
	Other	-	-	-

+P		Analyzer		
		Isch	Hrr	Other
Ref.	Isch	g	h	-
	Hrr	i	j	-
	Other	k	l	-

If considering both ischaemic and heart rate related ST change episodes together as **ST change episodes of unique type**, then the performance matrices can easily be reduced back to two-by-two, with:

$$\begin{aligned}
 \blacktriangleright \quad & TP_s = a+b+d+e & FN &= c+f \\
 & TP_p = g+h+i+j & FP &= k+l \\
 & Se = TP_s / (TP_s + FN) \\
 & +P = TP_p / (TP_p + FP)
 \end{aligned}$$

# Estimating performance of classifying ischaemic and non-ischaemic heart-rate related ST episodes

		Analyzer	
		Isch	Hrr
Ref.	Isch	TP	FN
	Hrr	FP	TN

		Analyzer	
		Isch	Hrr
Ref.	Isch	g	h
	Hrr	i	j

Analyzer-annotated ST segment episodes can further be classified between ischaemic and non-ischaemic heart rate related ST episodes. Classic performance measures in terms of  $Se$  and  $Sp$  can be used.

$$Se = \frac{TP}{TP + FN} = \frac{g}{g + h}$$

$$Sp = \frac{TN}{TN + FP} = \frac{j}{j + i}$$

# Estimating performance of classifying ischaemic and non-ischaemic heart-rate related ST episodes

- Analyzer using time-domain, KLT and Legendre Polynomial Transform approaches

Groups	Gross		Average	
	$Se[\%]$	$Sp[\%]$	$Se[\%]$	$Sp[\%]$
HR, ST	97.5	84.2	95.9	81.8
HR, ST, MD	97.8	85.0	98.4	82.5
HR, LPC	98.5	85.5	97.9	80.4
HR, LPC, MD	98.4	85.9	98.1	85.2

# Performance metrics

- The evaluation of an ST detection algorithm or analyzer should answer the following questions:
  - How well are ST episodes detected ?
  - How reliably is ST episode or ischaemic ST episode duration measured ?
  - How accurately are ST deviations measured ?
  - How well are ischaemic and non-ichaemic heart rate related ST episodes differentiated ?
  - **How well will the ST analyzer perform in the “real world” ?**

# Assessing the robustness (predicting the “real world” performance)

- Aggregate gross statistics
- Aggregate average statistics
- “Bootstrap method” of random generating new databases
- Noise stress test (assessing performance after adding noise to records)
- Sensitivity analysis by modifying analyzer's architecture parameters

# Assessing the robustness (predicting the “real world” performance)

- **Aggregate gross statistics**
  - How well an analyzer detects a randomly chosen ST episode ?
- **Aggregate average statistics**
  - How well the analyzer performs on a randomly chosen ambulatory record ?
- “Bootstrap method” of random generating new databases
- Noise stress test (assessing performance after adding noise to records)
- Sensitivity analysis by modifying analyzer's architecture parameters

# Assessing the robustness (predicting the “real world” performance)

<i>Technique</i>	↓	<i>ESC DB</i>				<i>LTST DB (Protocol B)</i>				
		<i>SE [%]</i>		<i>SD [%]</i>		<i>SE [%]</i>		<i>SD [%]</i>		
		<i>Se</i>	<i>+P</i>	<i>Se</i>	<i>+P</i>	<i>Se</i>	<i>+P</i>	<i>Se</i>	<i>+P</i>	
Time domain	[g]	81	76	-	-	-	-	-	-	
	[a]	84	81	-	-	-	-	-	-	
RMS method	[g]	-	-	-	-	-	-	-	-	
	[a]	84.7	86.1	75.3	68.2	-	-	-	-	
Time domain	[g]	79.2	81.4	-	-	-	-	-	-	
	[a]	81.5	82.5	-	-	-	-	-	-	
KLT approach	[g]	85.2	86.2	75.8	78.0	→	77.0	58.8	48.5	47.8
	[a]	87.1	87.7	78.2	74.1	→	74.0	61.4	54.8	58.4
→ Time domain, KLT	[g]	77.2	86.3	67.5	69.2	←	79.6	78.3	68.4	67.3
	[a]	81.3	89.2	77.6	68.9	←	78.9	80.7	73.1	74.9
Neural net	[g]	85.0	68.7	73.0	69.5		-	-	-	-
	[a]	88.6	78.4	72.2	67.5		-	-	-	-
Neural net, KLT	[g]	-	-	-	-		-	-	-	-
	[a]	77	86	-	-		-	-	-	-

# Assessing the robustness (predicting the “real world” performance)

- Aggregate gross statistics
- Aggregate average statistics
- “Bootstrap method” of random generating new databases
  - Is the analyzer's performance critically dependent on the choice of the test database ?
- Noise stress test (assessing performance after adding noise to records)
- Sensitivity analysis by modifying analyzer's architecture parameters



# Assessing the robustness (predicting the “real world” performance)

“**Bootstrap**” method (necessary assumption: the records of the origin(al) database are representative set of records for the problem domain ! ):

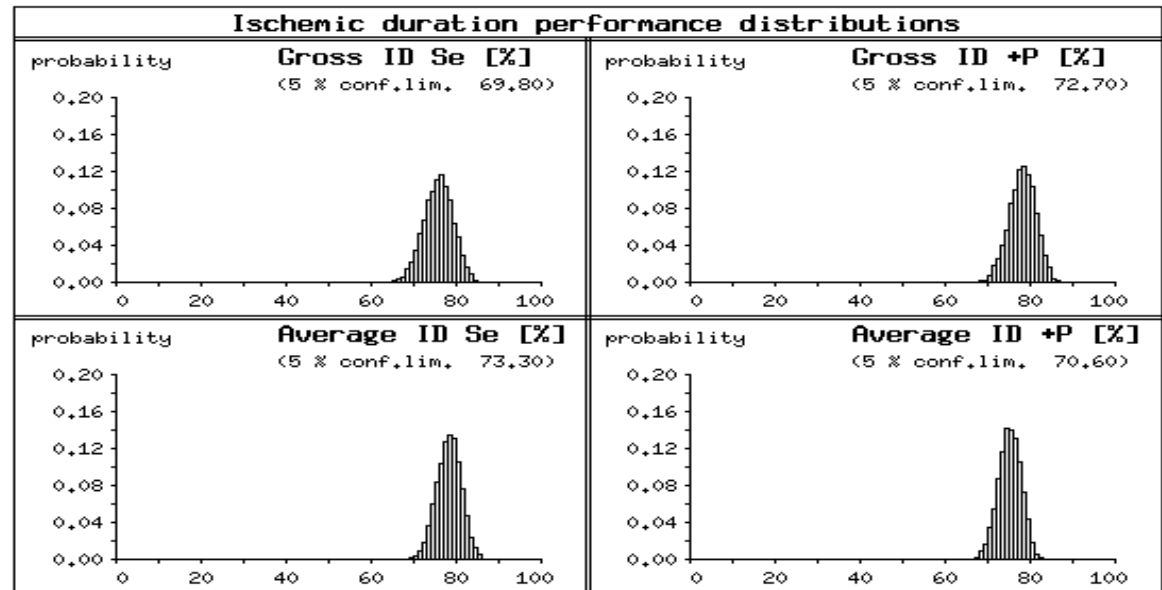
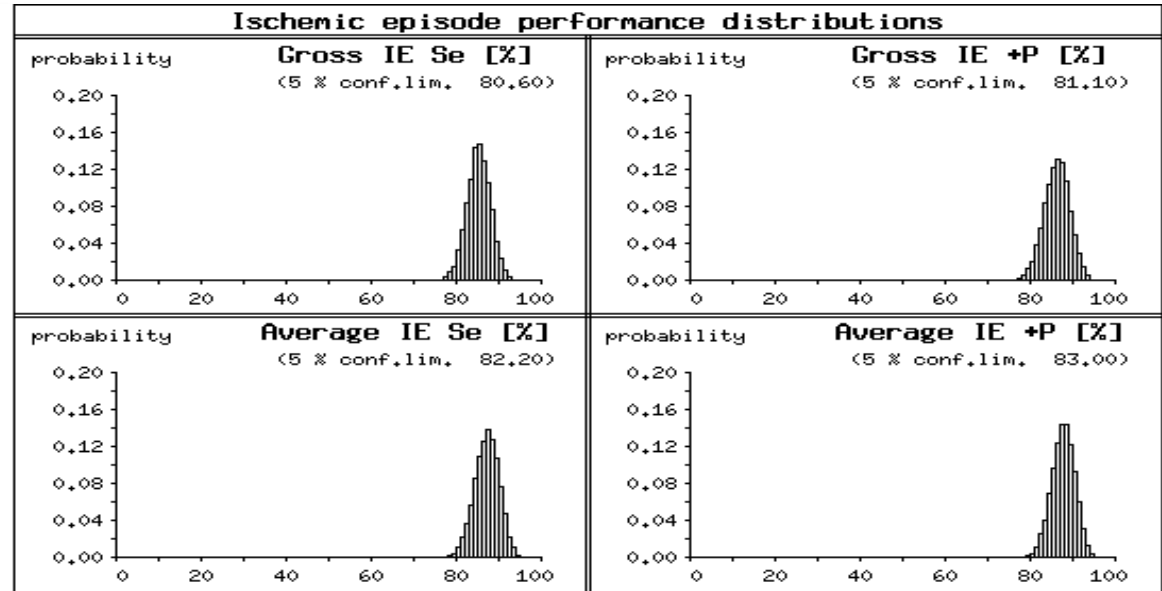
- (1) Randomly (and with replacement) choose  $L$  elements from the original database and insert them into a hypothetical database
- (2) Using the hypothetical database estimate the analyzer's performance
- (3) Repeat steps (1) and (2) many (10.000) times
- (4) Assess the estimates of the performance *distributions* obtained in the step 2 in the sense of 95% confidence limits

# Assessing the robustness (predicting the “real world” performance)

<i>Technique</i>		<i>ESC DB</i>				<i>LTST DB (Protocol B)</i>				
		<i>SE [%]</i>		<i>SD [%]</i>		<i>SE [%]</i>		<i>SD [%]</i>		
		<i>Se</i>	<i>+P</i>	<i>Se</i>	<i>+P</i>	<i>Se</i>	<i>+P</i>	<i>Se</i>	<i>+P</i>	
Time domain	[g]	81	76	-	-	-	-	-	-	
	[a]	84	81	-	-	-	-	-	-	
RMS method	[g]	-	-	-	-	-	-	-	-	
	[a]	84.7	86.1	75.3	68.2	-	-	-	-	
Time domain	[g]	79.2	81.4	-	-	-	-	-	-	
	[a]	81.5	82.5	-	-	-	-	-	-	
→ KLT approach	[g]	85.2	86.2	75.8	78.0	→	77.0	58.8	48.5	47.8
	[a]	87.1	87.7	78.2	74.1	→	74.0	61.4	54.8	58.4
Time domain, KLT	[g]	77.2	86.3	67.5	69.2	←	79.6	78.3	68.4	67.3
	[a]	81.3	89.2	77.6	68.9	←	78.9	80.7	73.1	74.9
Neural net	[g]	85.0	68.7	73.0	69.5	-	-	-	-	
	[a]	88.6	78.4	72.2	67.5	-	-	-	-	
Neural net, KLT	[g]	-	-	-	-	-	-	-	-	
	[a]	77	86	-	-	-	-	-	-	

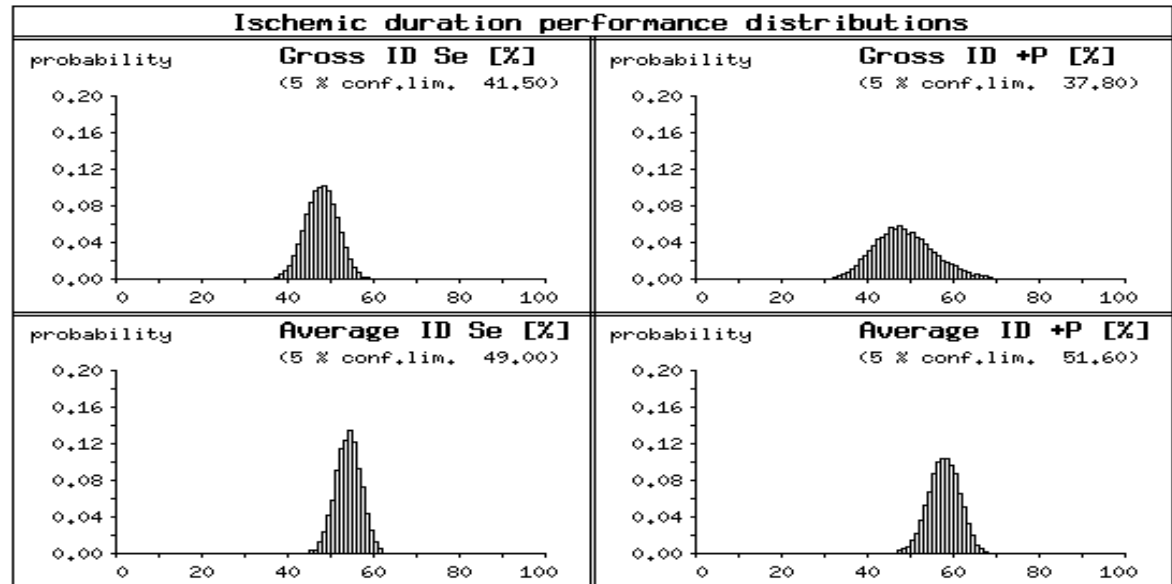
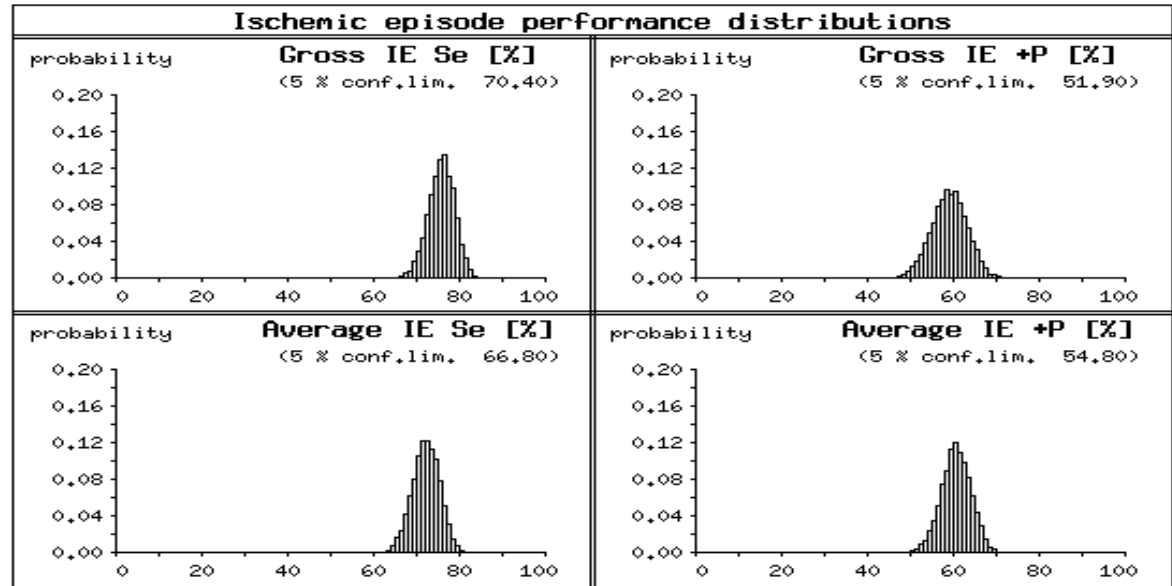
# Assessing the robustness (predicting the “real world” performance)

- **KLT approach**, “bootstrap” distributions as obtained on ESC DB (development database)



# Assessing the robustness (predicting the “real world” performance)

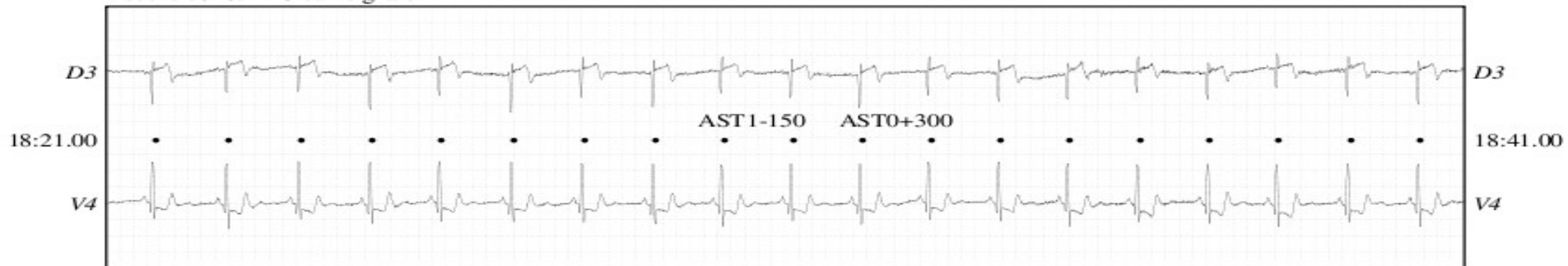
- **KLT approach**, “bootstrap” distributions as obtained on LTST DB (testing database)



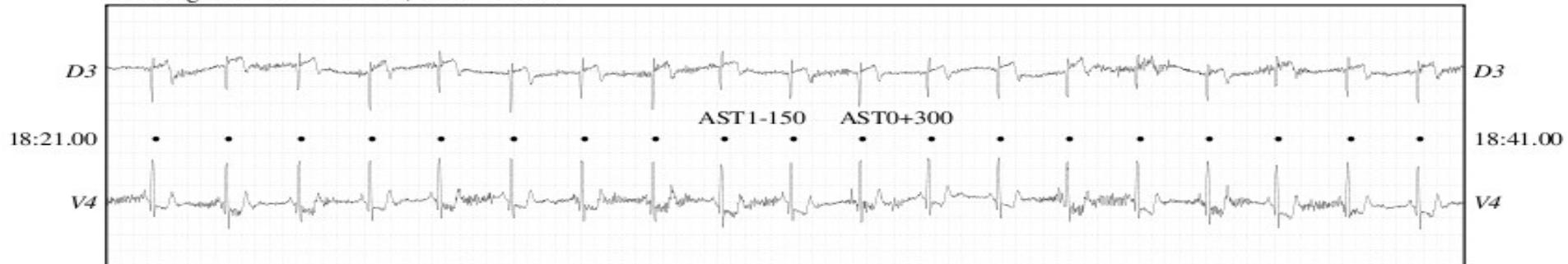
# Assessing the robustness (predicting the “real world” performance)

- Aggregate gross statistics
- Aggregate average statistics
- “Bootstrap method” of random generating new databases
- Noise stress test (assessing performance after adding noise to records)
  - What is the minimum critical signal-to-noise ratio at which the analyzer's performance is still acceptable ?
- Sensitivity analysis by modifying analyzer's architecture parameters

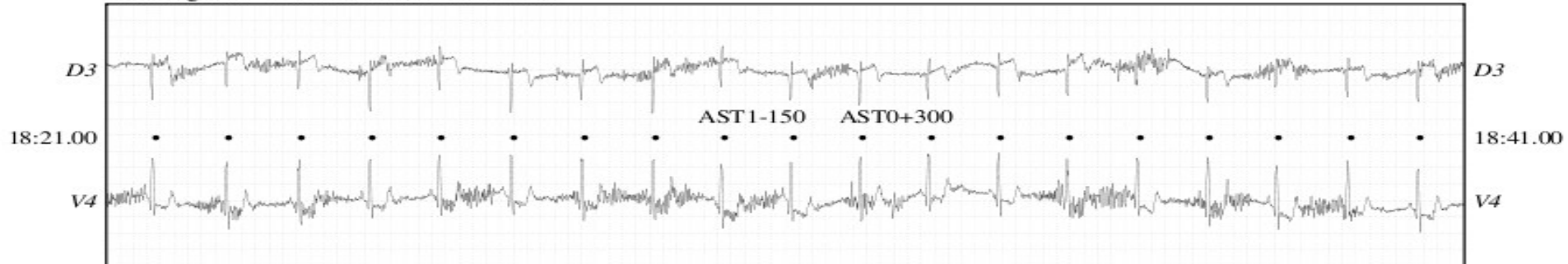
Record e0107 "Clean signals"



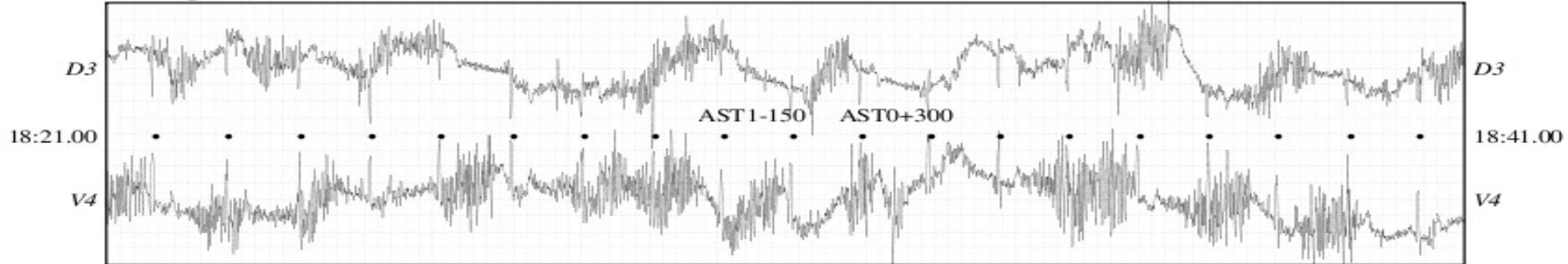
Record g0107 Muscle noise, SNR = 24 dB



Record g0107 Muscle noise, SNR = 18 dB



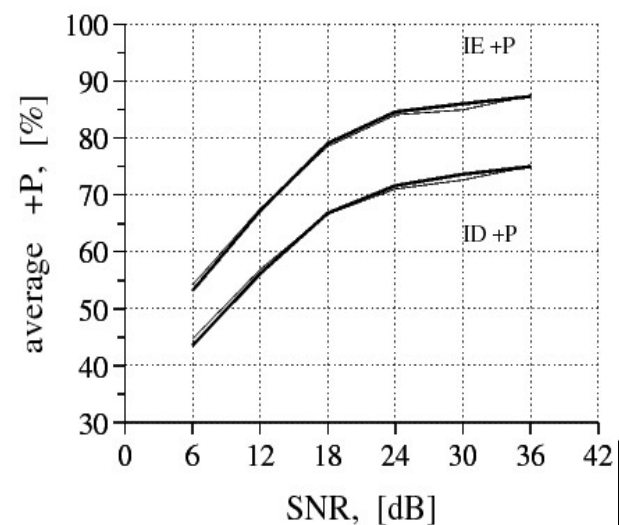
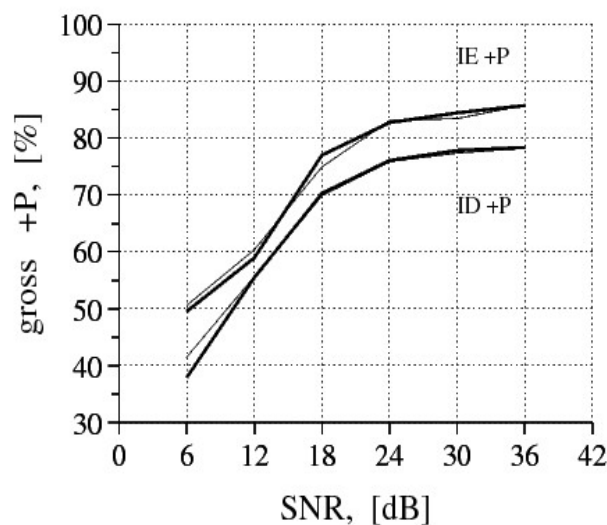
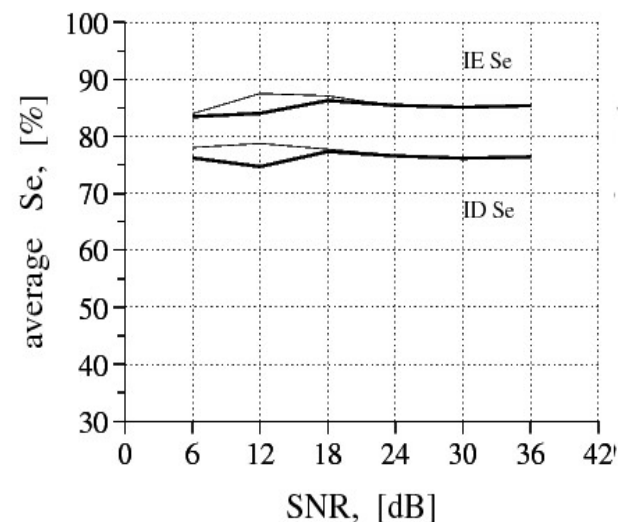
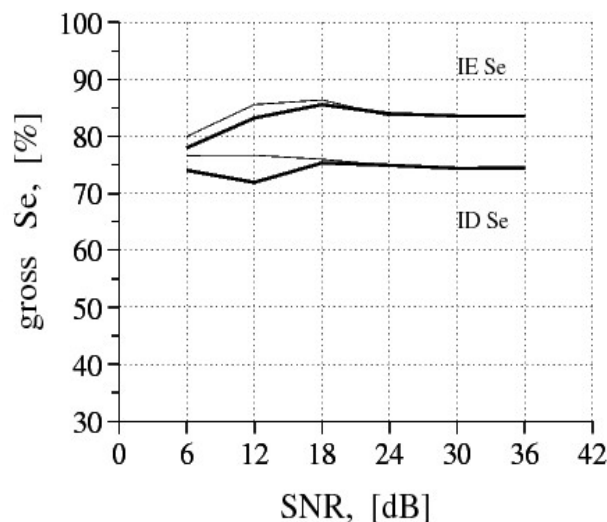
Record g0107 Muscle noise, SNR = 6 dB



# Assessing the robustness (predicting the “real world” performance)

- **KLT approach**, “bootstrap” distributions as obtained on ESC DB (development database)

Influence of noise stress test



# Assessing the robustness (predicting the “real world” performance)

- Aggregate gross statistics
- Aggregate average statistics
- “Bootstrap method” of random generating new databases
- Noise stress test (assessing performance after adding noise to records)
- **Sensitivity analysis by modifying analyzer's architecture parameters**
  - Are the architecture parameters critically tuned to the development database ?



# Assessing the robustness (predicting the “real world” performance)

- **KLT approach**, “bootstrap” distributions as obtained on ESC DB (development database)

Influence of KLT feature-vector dimensionality

