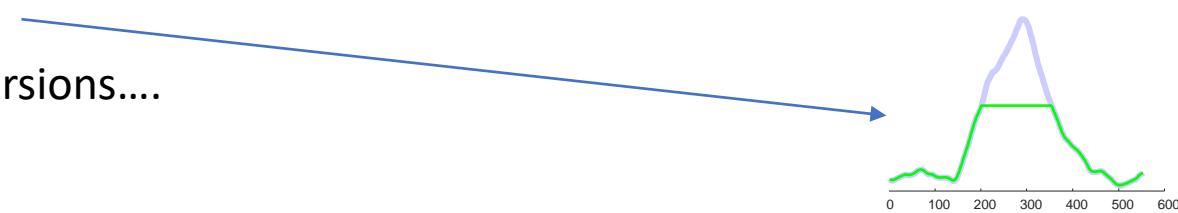


# HEX UCR Anomaly Benchmark Datasets 2021

# Dataset Design Principles I

- **Remove the threshold question.** Anomaly detection typically consists of two parts. A) Find the candidate region(s) that *might be* anomalies. B) Test if, under your model, you should flag these as anomalies. Here we remove ‘B’, by telling the world that there is *exactly one anomaly* in the test data. We do this because ‘B’ can depend on external factors (misclassification costs etc.), and we think that in most domains, if you can do ‘A’ robustly, ‘B’ will be easy.
- **Try to have diverse datasets.**
- **Use real data only in the case** you can be sure the anomaly is the *only* (or by a large margin, *most significant*) anomaly.
- **For synthetic data, model something in the real world** (when possible). For example:
  - This anomaly models the actor falling down.
  - This anomaly models a terrorist attack in Melbourne, and the police blocking the road for an hour.
  - This anomaly models a nurse placing her hand under the respiration strap.
- **Avoid Goldilocks.** For at least some problems, make multiple versions....
  - One that is obvious, probably *any* algorithm can find them.
  - One that is more subtle
  - One that is very very subtle, probably *no* algorithm can find them.
- **Implication of the above:** It is virtually certain that 100% accuracy is not possible.

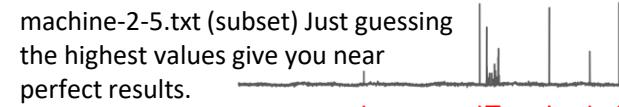


This anomaly models a nurse placing her hand under the respiration strap

UCR\_Anomaly\_respiration1\_100000\_110260\_110412.txt

# Dataset Design Principles II

- **Only one anomaly per dataset.** Some datasets (see machine-2-5.txt) have so many anomalies that it confuses scoring metrics.
- **Anomalies can appear anywhere.** Some datasets (esp. Yahoo) have most of the anomalies at the end (run-to-failure). But this is an unfair clue to give algorithms.
- **Anomalies should not (in general) be the highest/lowest values.** Many datasets (Yahoo/ServerMachineDataset/NASA) have some problems that are trivial to solve with the “pick the highest value” heuristic.
- **(More general version of previous rule)** Anomalies should not (in general) have very different means/standard deviations/skewness or other simple summary statistics. Such things *can* be anomalies, but we can typically find them with 60-year old change detection algorithms.
- **Avoid always inserting anomalies at round numbers,** say at location 10,000 or 11,500. As that gives a clue to the algorithms.



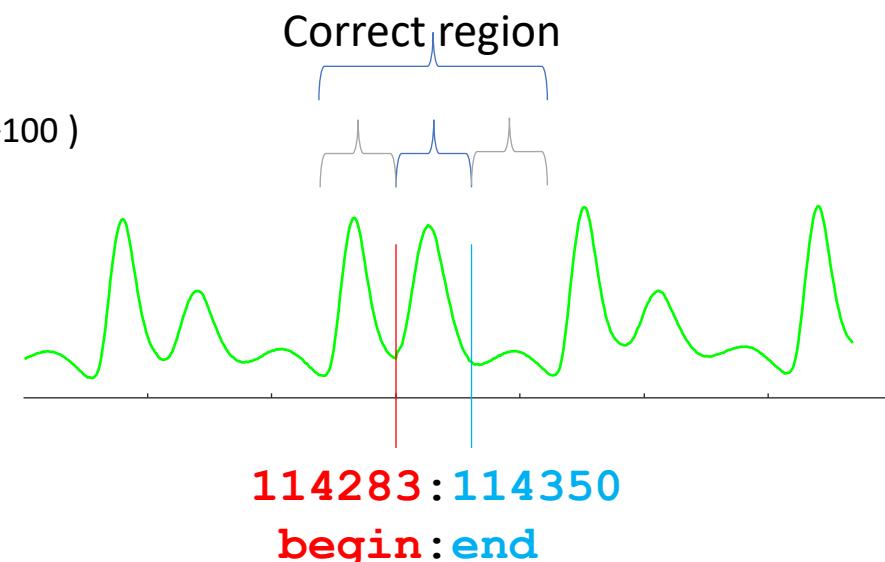
# Scoring Function Design Principles I

**Avoid complex and opaque scoring functions** We want a scoring function that..

- Is a single number, for easy comparisons.
- Does not have spurious location precision. If the ground truth says the anomaly is at say 1250, and an algorithm reports 1247 or 1254, it should be counted as correct. This problem is compounded by the fact that different algorithms report the *leading edge*, the *center* or the *trailing edge* of a sliding window.
- Has a binary score for each example, that can be combined to a real number for the full collection.
- Reports a number close to zero for a “random dart” algorithm (i.e. the default rate) and close to one for a perfect algorithm.

My suggestion

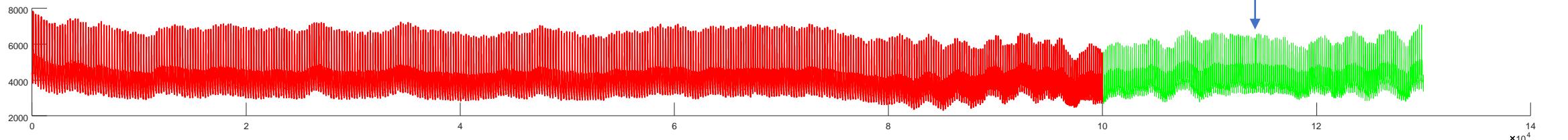
- Let length of anomaly be  $L$ , 
$$L = \text{end} - \text{begin} + 1$$
- Let the prediction of an algorithm be an integer  $P$
- $P$  is labeled as correct if: 
$$\min(\text{begin}-L, \text{begin}-100) < P < \max(\text{end}+L, \text{end}+100)$$
- Why the ‘100’ case? Some anomalies can be as short as a single point.



- Dataset number
- Mnemonic name
- From 1 to X is training data
- Begin anomaly
- End anomaly

The data comes from a healthy male on a tilt table.  
At first, he is supine, at around 80000, the table is tilted forward. The trace is his APB.

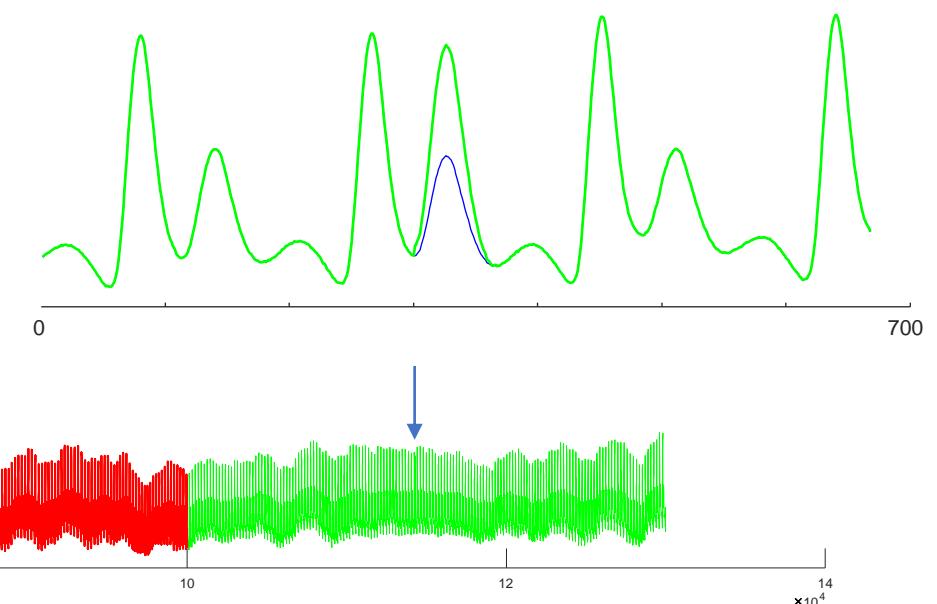
The anomaly is synthetic. There is a secondary peak after the dicrotic notch. It is normally about half the size of the peak systolic pressure. For one randomly chosen beat, we made it much greater, almost as big as the main peak.



## Sample format

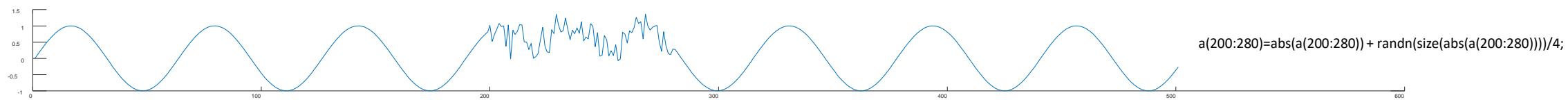
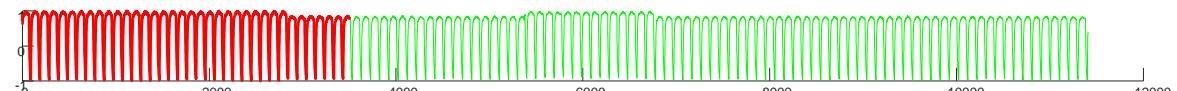
Note the structure of the file names  
Files are '-ascii' format.  
Note the layout of the plots

Blue is original data; green is data after anomaly was introduced



# Notes

- Every dataset should be considered completely independent of the other datasets. Any algorithm that considers data from two or more datasets at the same time should be considered cheating.
- However, an algorithm *could* use external data. For example
  - An algorithm could “realize” the dataset is an ECG (without human intervention, obviously), it could then search the web for similar ECGs, and use that information to help build the model.
  - Note, we do not think this is likely to help, but it is within the spirit of fair evaluation
- You can assume:
  - The **training data** is free of anomalies
  - The test data has only one anomaly. Or equivalently, if it has more than one, exactly one is much much more obvious/significant than the others.
- You should **predict the center** of the anomaly:
  - Some algorithms, by default, output the beginning, the middle, or the end of the anomaly. A good scoring function will have some “wiggle room” to allow for this, so this is probably a non-issue.
  - However, it is probably best to output the *center* of where your algorithm thinks the anomaly is. For example, in the below, the anomaly starts at 200, and is of length 80. It would be best to output 240.



# Known issues

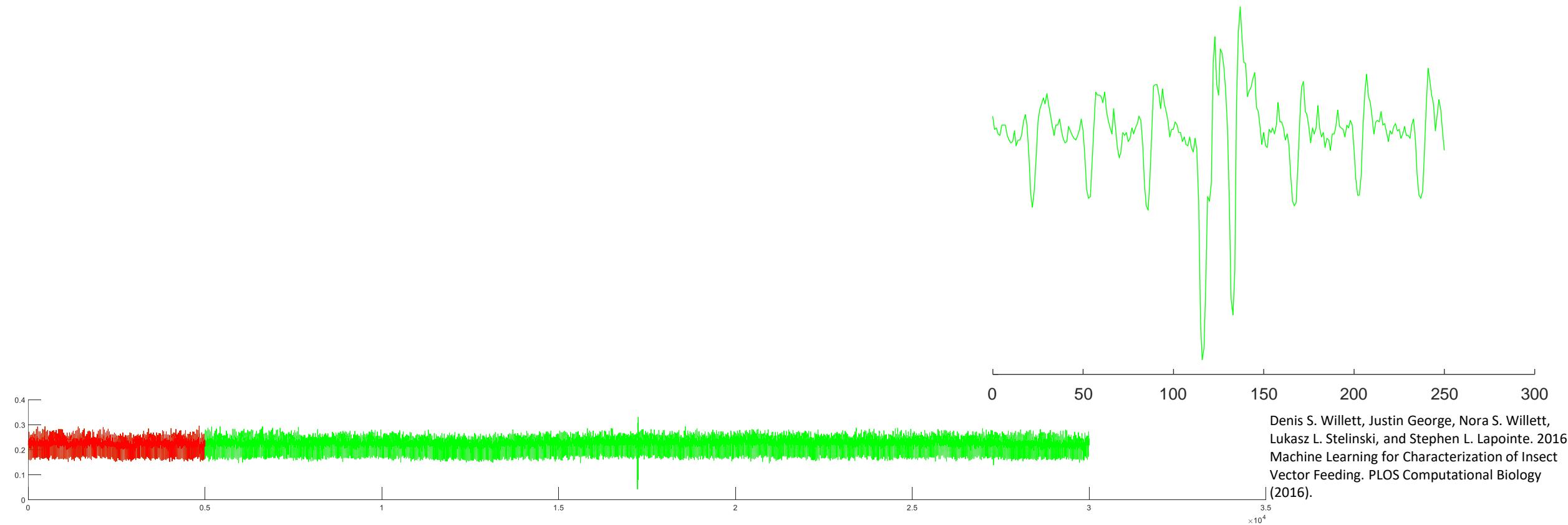
- This will be updated if new issues are brought to our attention
- Some files are formatted differently: tabulator between values instead of line break. This can lead to a problem when reading the files, e.g. when using `pandas.read_csv(...)`. This applies to file number 204, 205, 206, 207, 208, 225, 226, 242 and 243.

# Key to Datasets

# UCR\_Anomaly\_Lab2Cmac011215EPG1\_5000\_17210\_17260.txt

This is a real dataset from an insect, and Asian Citrus Psyllid, recorded using an EPG apparatus.

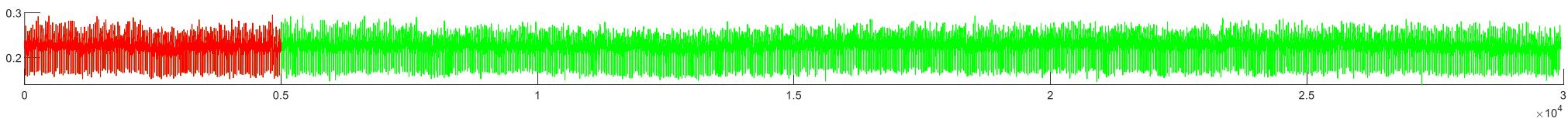
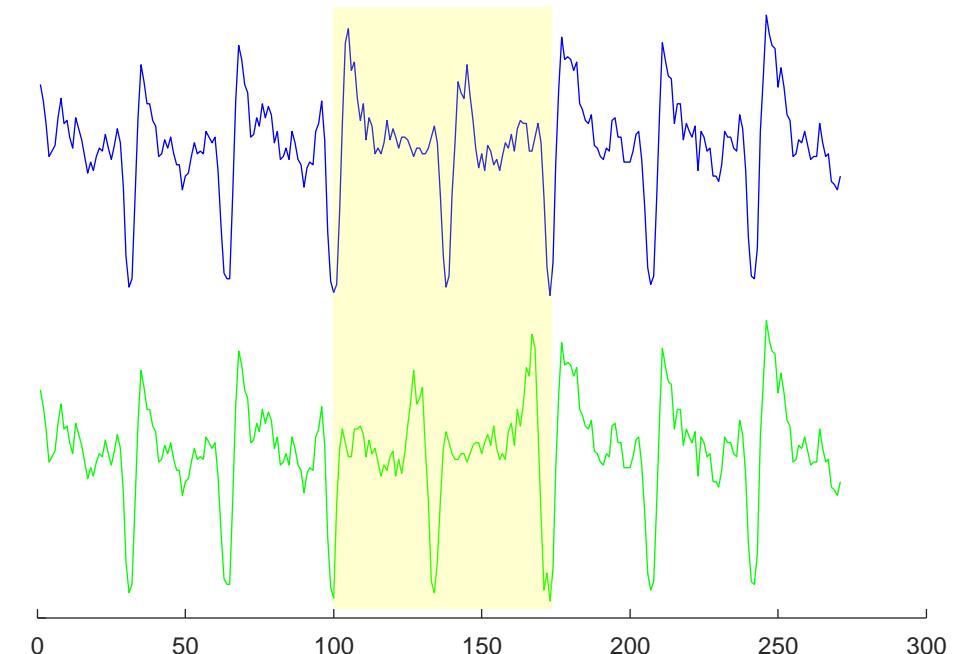
Here the anomaly is completely natural, probably the insect moved its stylet to a new vein.



# UCR\_Anomaly\_Lab2Cmac011215EPG2\_5000\_27862\_27932.txt

This is a real dataset from an insect, and Asian Citrus Psyllid, recorded using an EPG apparatus.

Here we removed the natural anomaly and then reversed the direction of two beats.

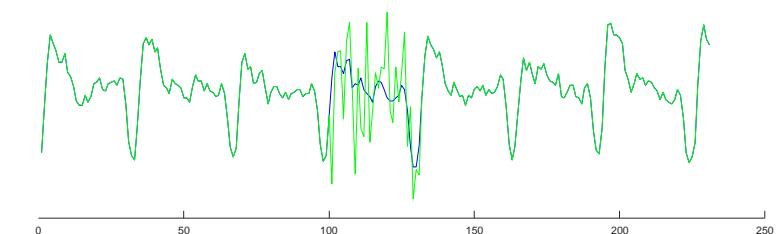
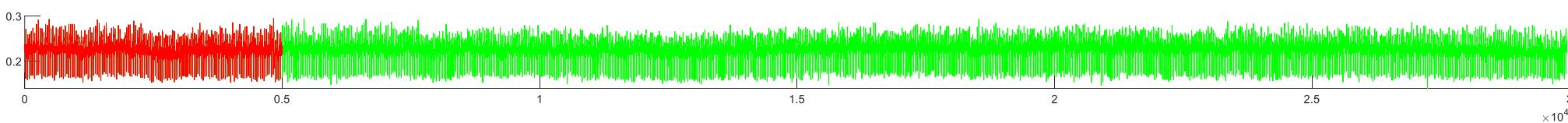


Denis S. Willett, Justin George, Nora S. Willett,  
Lukasz L. Stelinski, and Stephen L. Lapointe. 2016.  
Machine Learning for Characterization of Insect  
Vector Feeding. PLOS Computational Biology  
(2016).

# UCR\_Anomaly\_Lab2Cmac011215EPG3\_5000\_16390\_16420.txt

This is a real dataset from an insect, and Asian Citrus Psyllid, recorded using an EPG apparatus.

Here we removed the natural anomaly and added a lot of noise to a single beat.



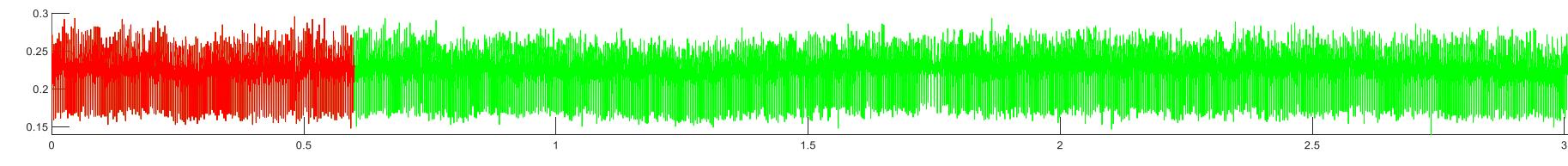
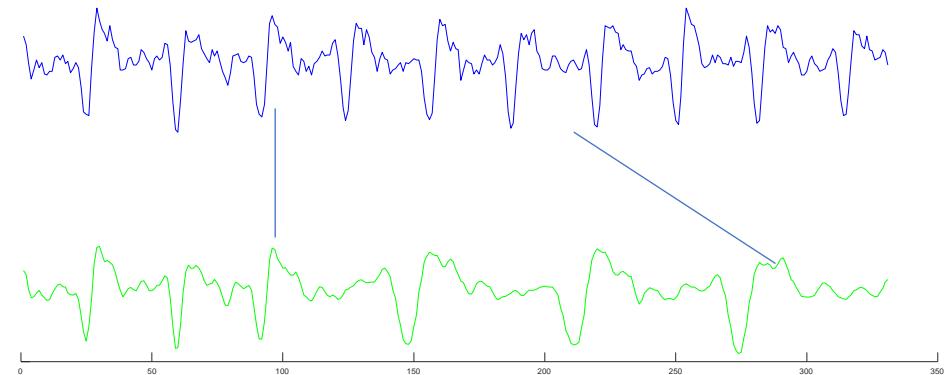
Denis S. Willett, Justin George, Nora S. Willett, Lukasz L. Stelinski, and Stephen L. Lapointe. 2016. Machine Learning for Characterization of Insect Vector Feeding. PLOS Computational Biology (2016).

# UCR\_Anomaly\_Lab2Cmac011215EPG4\_6000\_17390\_17520.txt

This is a real dataset from an insect, and Asian Citrus Psyllid, recorded using an EPG apparatus.

Here we removed the natural anomaly and then “slowed down” about three beats with interpolation.

We also truncate a few random points from the beginning, just so this would not align with the other datasets in this series.



Denis S. Willett, Justin George, Nora S. Willett, Lukasz L. Stelinski, and Stephen L. Lapointe. 2016. Machine Learning for Characterization of Insect Vector Feeding. PLOS Computational Biology (2016).

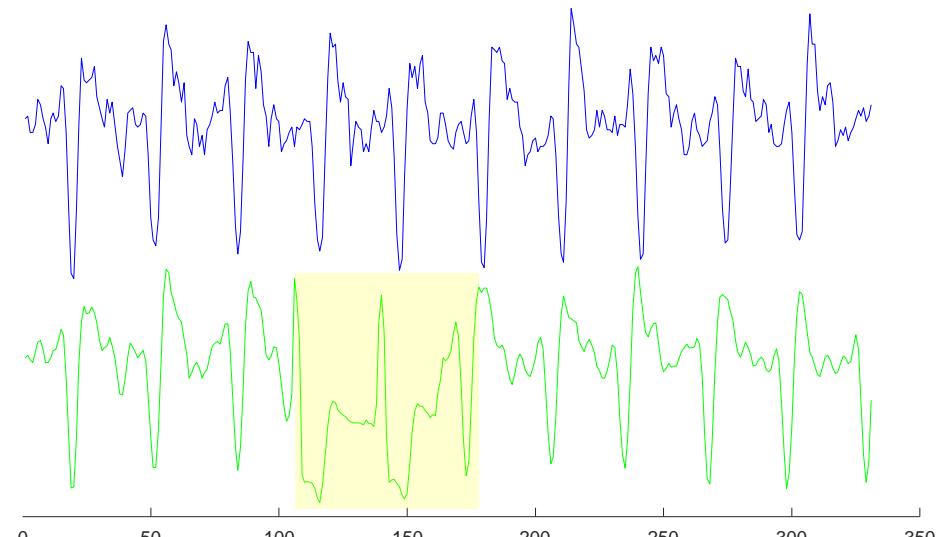
# UCR\_Anomaly\_Lab2Cmac011215EPG5\_7000\_17390\_17520.txt

This is a real dataset from an insect, and Asian Citrus Psyllid, recorded using an EPG apparatus.  
Here we removed the natural anomaly.

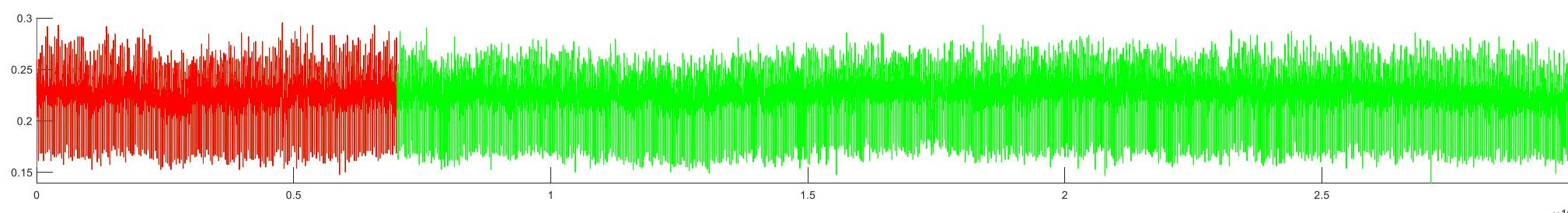
We inserted about 1.5 human heartbeats (from UCR\_Anomaly\_ECG3\_8000\_17000\_17100.txt)

We made a quick effort to edit the human beats to have approximately the same mean, variance and periodicity as the original data, and to smooth the joins.

We also truncate a few random points from the beginning, just so this would not align with the other datasets in this series.



0 50 100 150 200 250 300 350



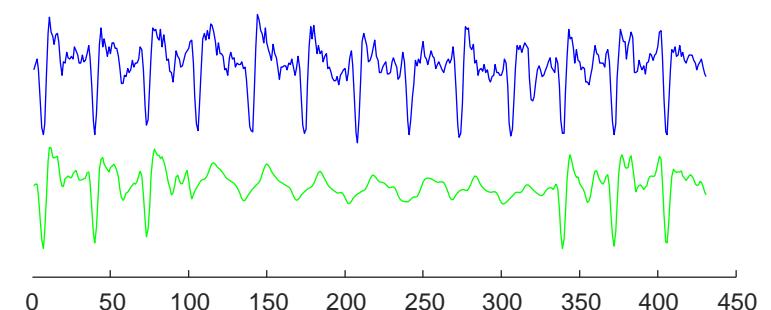
Denis S. Willett, Justin George, Nora S. Willett, Lukasz L. Stelinski, and Stephen L. Lapointe. 2016. Machine Learning for Characterization of Insect Vector Feeding. PLOS Computational Biology (2016).

# UCR\_Anomaly\_Lab2Cmac011215EPG6\_7000\_12190\_12420.txt

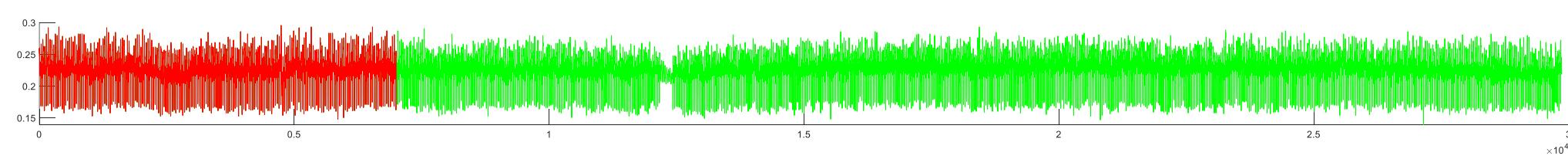
This is a real dataset from an insect, and Asian Citrus Psyllid, recorded using an EPG apparatus.  
Here we removed the natural anomaly.

We greatly smoothed a section. This both smooths the data, and reduces the variance.

We also truncate a few random points from the beginning, just so this would not align with the other datasets in this series.



0 50 100 150 200 250 300 350 400 450



Denis S. Willett, Justin George, Nora S. Willett, Lukasz L. Stelinski, and Stephen L. Lapointe. 2016. Machine Learning for Characterization of Insect Vector Feeding. PLOS Computational Biology (2016).

# UCR\_Anomaly\_TkeepFirstMARS\_3500\_5365\_5380.txt

This is a real dataset from NASA spacecraft, that appeared in a KDD 2018 paper.

We joined the train and test sets.

The original dataset had two anomalies.

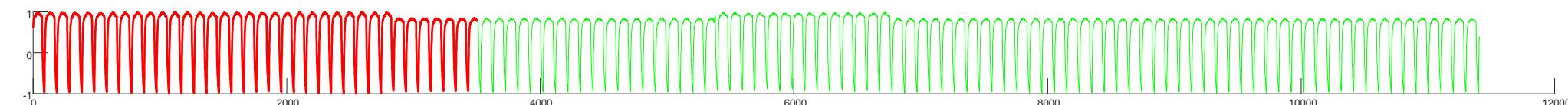
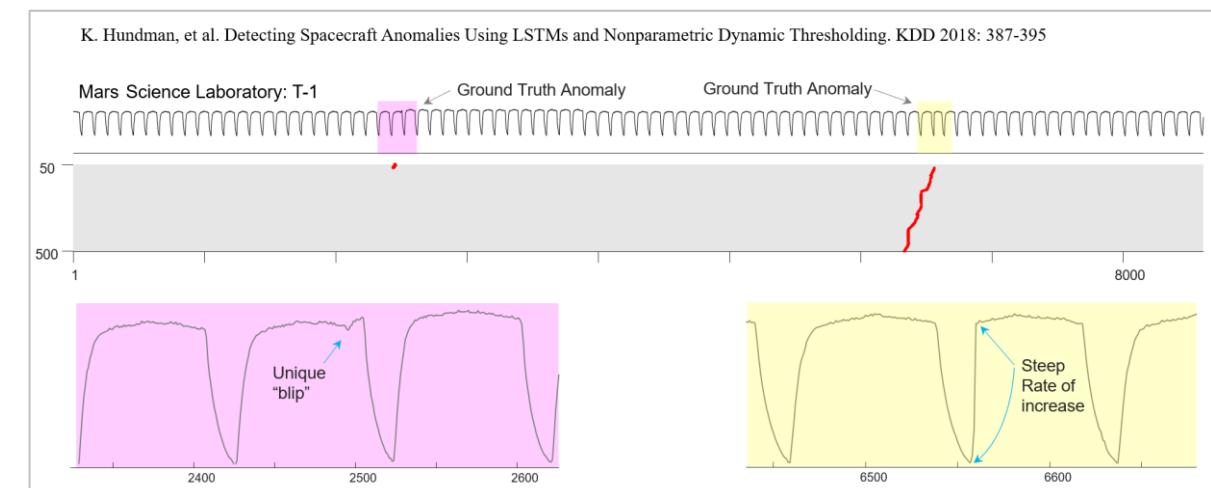
We carefully edited the data, into two datasets, each with *one* anomaly.

This one contains “**unique blip**”.

There is also a small amount of level change in this datasets.

However, note that it occurs in both the train and test sets, so it is not the anomaly.

This figure is from the original test split



# UCR\_Anomaly\_TkeepSecondMARS\_3500\_9330\_9340.txt

This is a real dataset from NASA spacecraft, that appeared in a KDD 2018 paper.

We joined the train and test sets.

The original dataset had two anomalies.

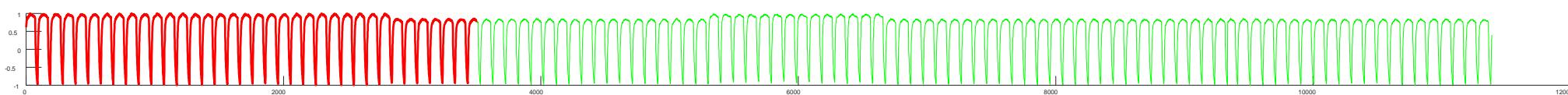
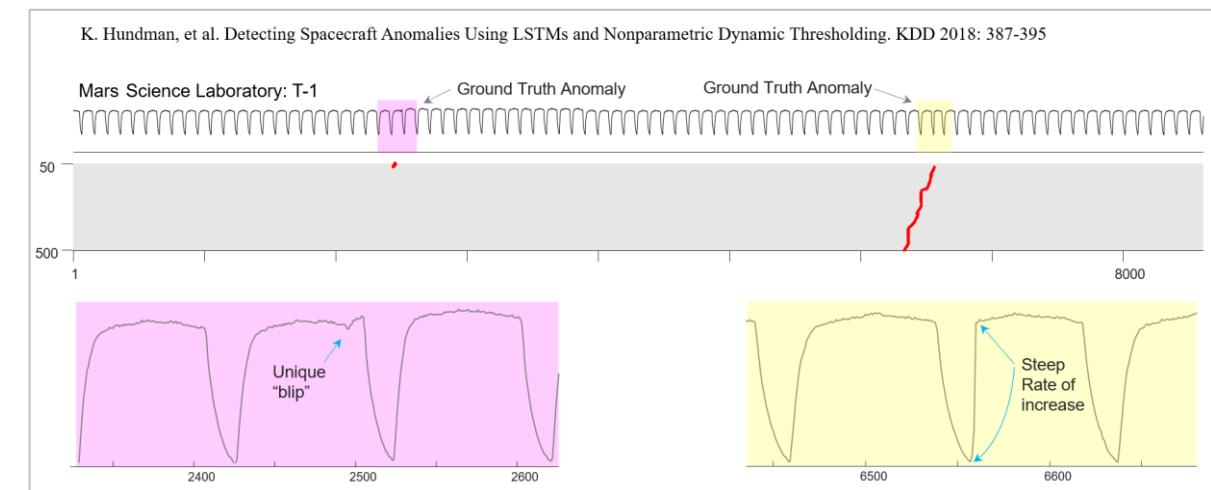
We carefully edited the data, into two datasets, each with *one* anomaly.

This one contains “[steep rate of increase](#)”.

There is also a small amount of level change in this datasets.

However, note that it occurs in both the train and test sets, so it is not the anomaly.

This figure is from the original test split



# UCR\_Anomaly\_TkeepThirdMARS\_3500\_4711\_4809.txt

This is a real dataset from NASA spacecraft, that appeared in a KDD 2018 paper.

We joined the train and test sets.

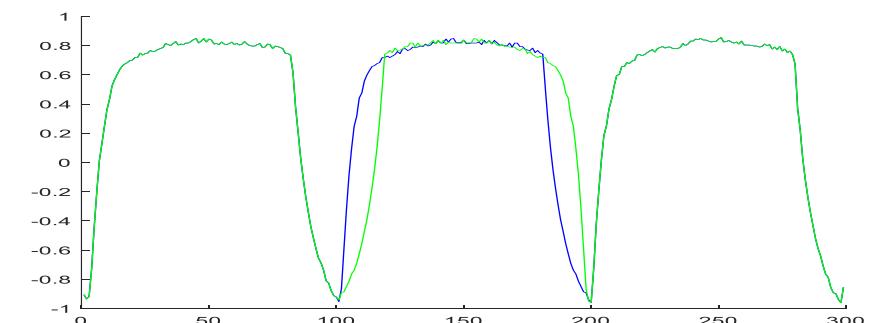
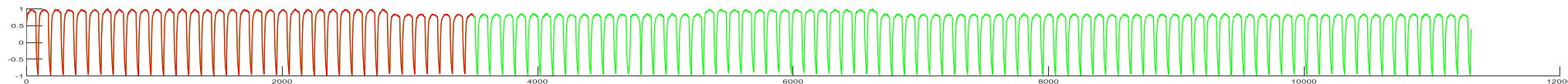
The original dataset had two anomalies.

We carefully removed both the anomalies, and added one of our own.

We flipped one cycle backwards

There is also a small amount of level change in this datasets.

However, note that it occurs in both the train and test sets, so it is not the anomaly.



# UCR\_Anomaly\_TkeepForthMARS\_3500\_5988\_6085.txt

This is a real dataset from NASA spacecraft, that appeared in a KDD 2018 paper.

We joined the train and test sets.

The original dataset had two anomalies.

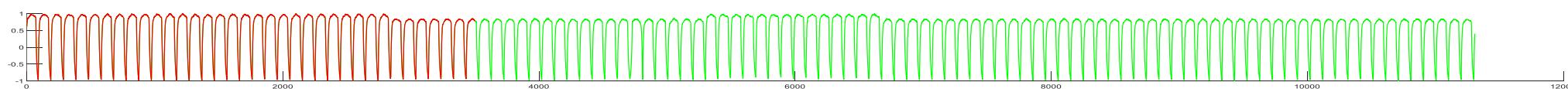
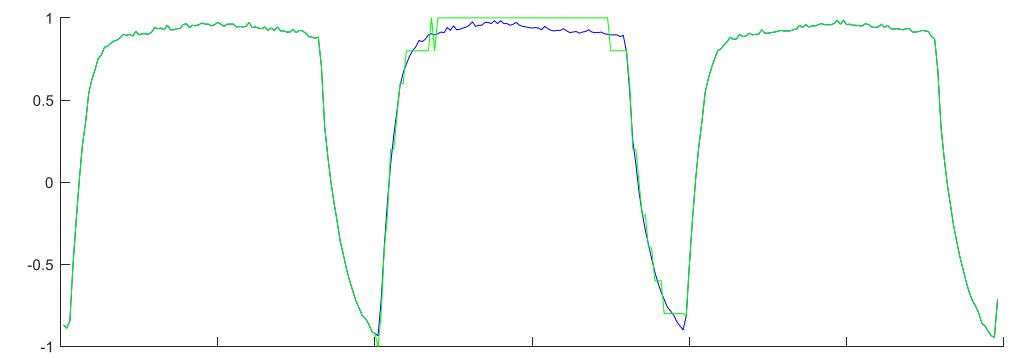
We carefully removed the anomalies, and added one of our own.

We reduced the bit depth of one cycle.

There is also a small amount of level change in this datasets.

However, note that it occurs in both the train and test sets, so it is not the anomaly.

```
sub = T(start_anomaly:end_anomaly) ;
sub = sub*5;
sub = round(sub);
sub = sub /5;
T(start_anomaly:end_anomaly) = sub;
```



# UCR\_Anomaly\_TkeepFifthMARS\_3500\_5988\_6085.txt

This is a real dataset from NASA spacecraft, that appeared in a KDD 2018 paper.

We joined the train and test sets.

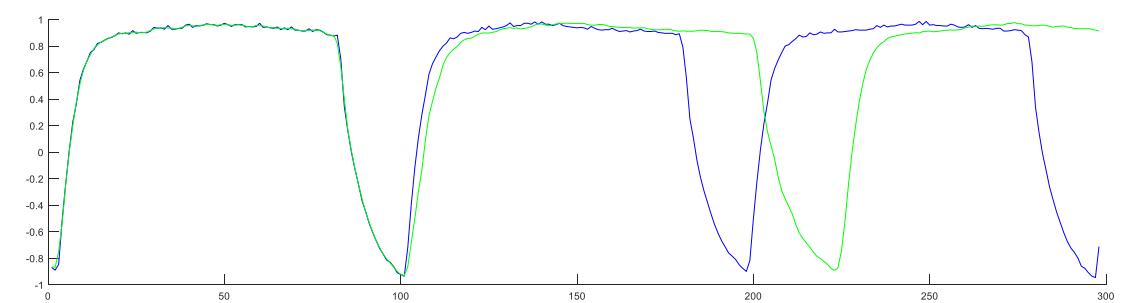
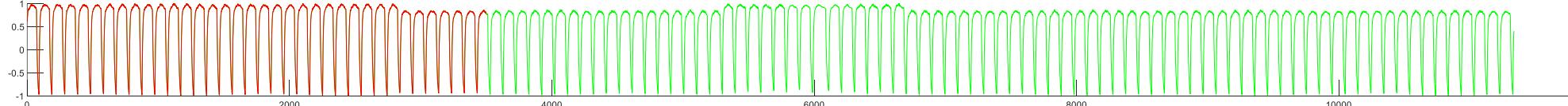
The original dataset had two anomalies.

We carefully removed the anomalies, and added one of our own.

We made one beat about 25% slower than the rest

There is also a small amount of level change in this datasets.

However, note that it occurs in both the train and test sets, so it is not the anomaly.



# UCR\_Anomaly\_gaitHunt1\_18500\_33070\_33180.txt

Selected  
input: record  
gaitndd/hunt17,  
from 0:00.000 to  
5:00.000

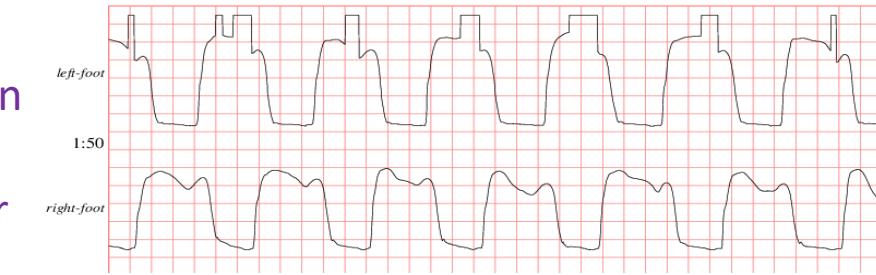
Gait in  
Neurodegenerative  
Disease Database  
(gaitndd)

This dataset comes from someone walking on a force plate in a biomechanics lab.

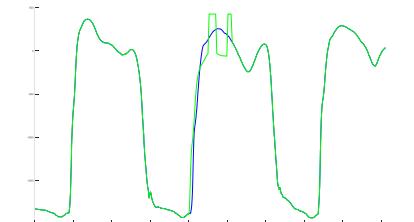
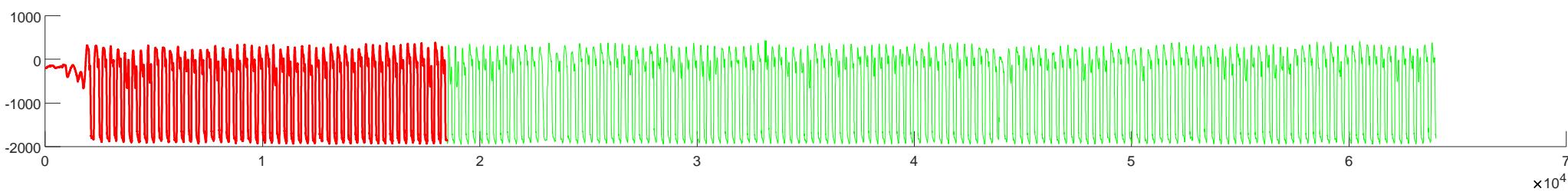
The individual had a mostly symmetric gait, however, after some time, the left foot sensor developed a fault.

This allows us to create an almost 100% natural dataset. We simply took faulty data from the left foot, and used it to replace some right foot data. We shifted it by a half cycle, slightly smoothing it to remove cut and paste artifacts, and reduced its amplitude so it was not necessarily the tallest peak.

Fault in  
left  
sensor



Blue is original data,  
green is data after  
anomaly was  
introduced

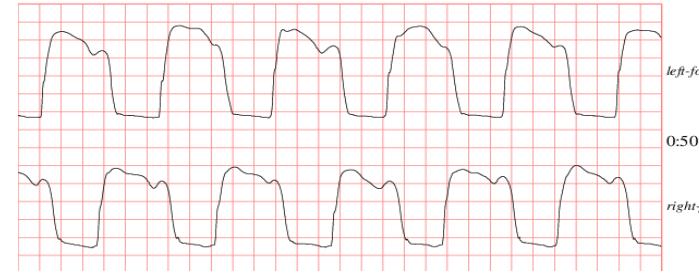


This dataset comes from someone walking on a force plate in a biomechanics lab.

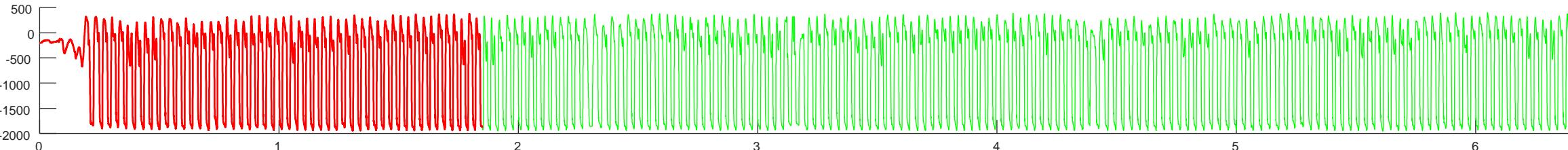
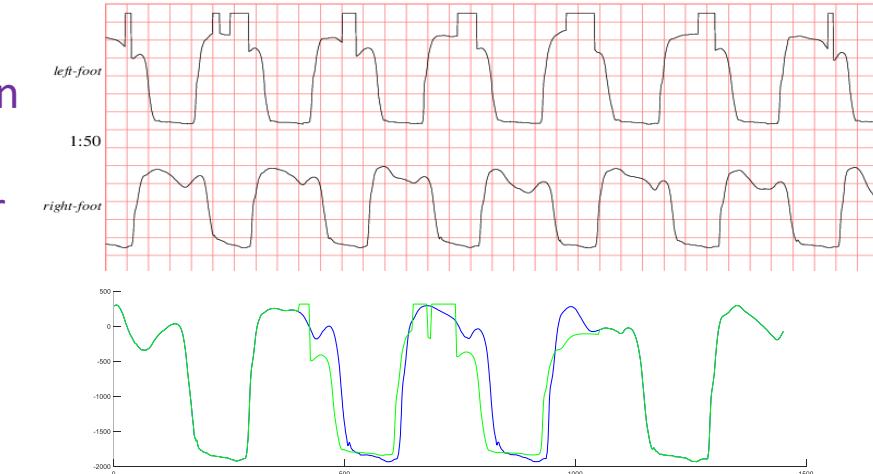
The individual had a mostly symmetric gait, however, after some time, the left foot sensor developed a fault.

This allows us to create an almost 100% natural dataset. We simply took faulty data from the left foot, and used it to replace some right foot data. We shifted it by a half cycle, slightly smoothing it to remove cut and paste artifacts, and reduced its amplitude so it was not necessarily the tallest peak.

First  
minute  
or so...



Fault in  
left  
sensor



# UCR\_Anomaly\_gaitHunt3\_23400\_38400\_39200.txt

Selected  
input: record  
gaitndd/hunt17,  
from 0:00.000 to  
5:00.000

Gait in  
Neurodegenerative  
Disease Database  
(gaitndd)

This dataset comes from someone walking on a force plate in a biomechanics lab.

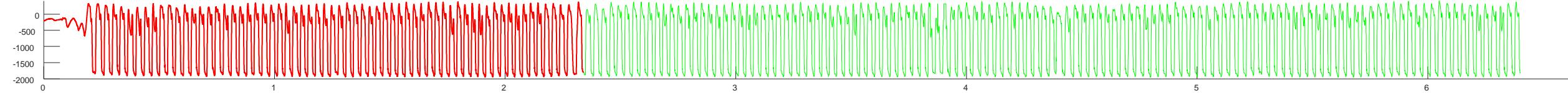
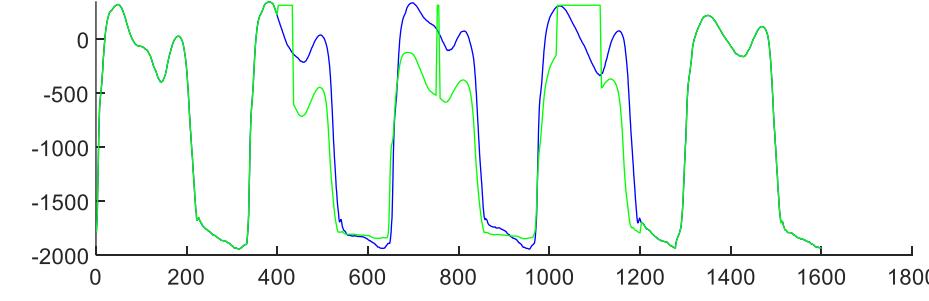
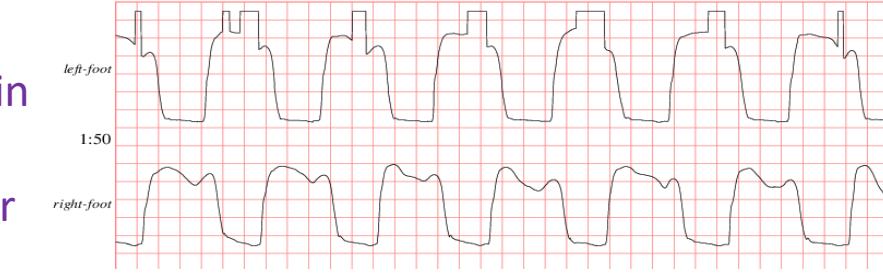
The individual had a mostly symmetric gait, however, after some time, the left foot sensor developed a fault.

This allows us to create an almost 100% natural dataset. We simply took faulty data from the left foot, and used it to replace some right foot data. We shifted it by a half cycle, slightly smoothing it to remove cut and paste artifacts, and reduced its amplitude so it was not necessarily the tallest peak.

First  
minute  
or so...



Fault in  
left  
sensor



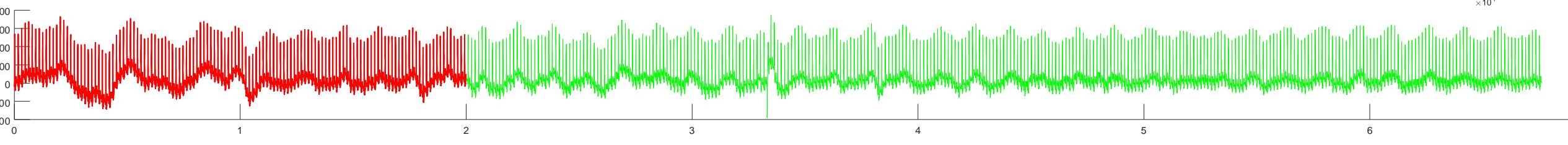
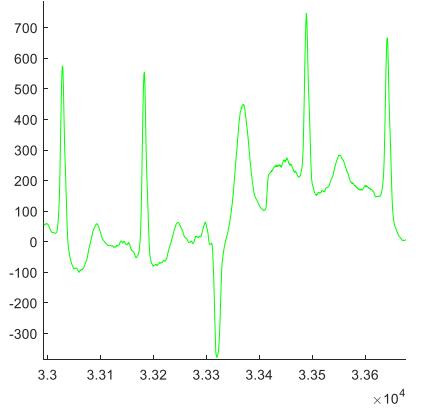
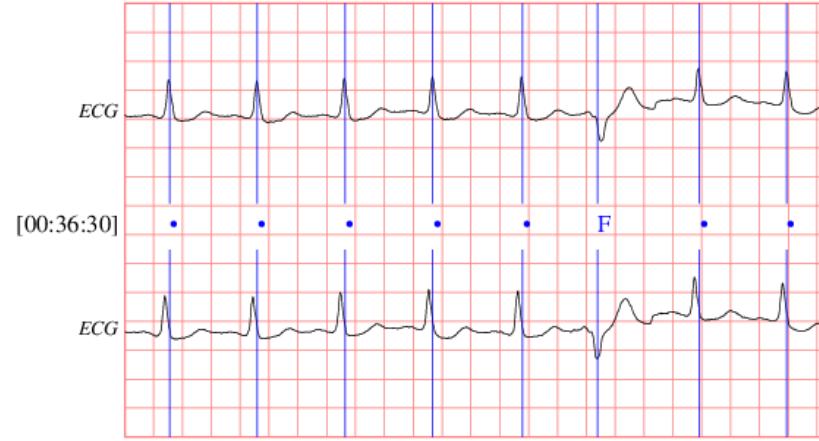
# UCR\_Anomaly\_sddbECG\_20000\_33270\_33400.txt

Selected  
input: record  
sddb/30 ,  
annotator atr ,  
from  
[00:34:20.000 1]  
to [01:34:20.000  
1]

[Sudden Cardiac  
Death Holter  
Database \(sddb\)](#)

This is a completely natural dataset. We found a long stretch of ECG that had a single anomaly (F = Fusion of ventricular and normal beat).

The data was *wandering baseline*, which I suspect will confuse many algorithms.



# UCR\_Anomaly\_gait1\_20000\_38500\_38800.txt

Selected  
input: record  
gaitndd/hunt14 ,  
from 0:00.000 to  
5:00.000

Gait in  
Neurodegenerativ  
e Disease  
Database  
(gaitndd)

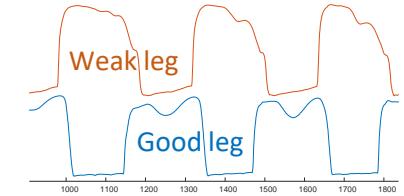
This dataset comes from someone walking on a force plate in a biomechanics lab.

The individual had Huntington's disease, and a highly asymmetric gait (a limp).

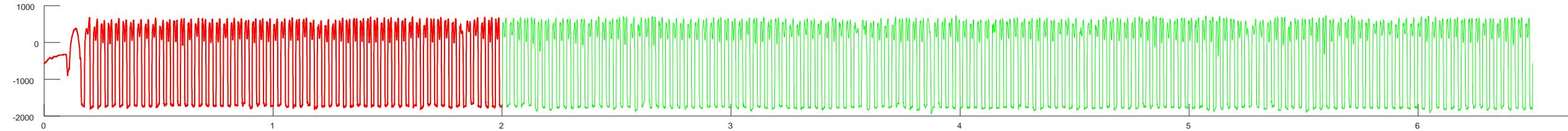
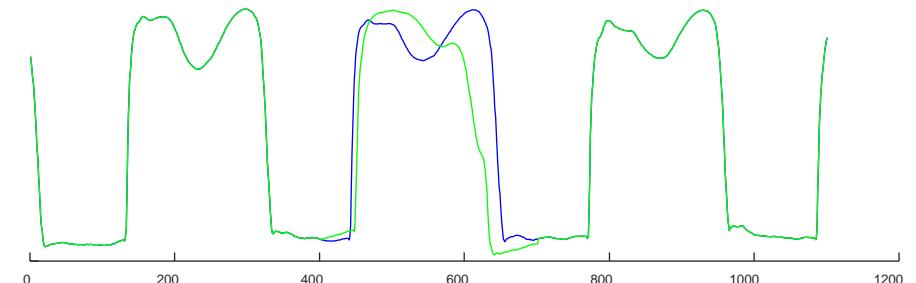
The dataset is from the good leg (his right)

The anomaly is that at a random location we swapped in the weak leg (shifting it by a half cycle, so it lines up correctly) for one cycle.

This models an individual having a sudden, but thankfully short lived, weakness or pain in their leg.



Blue is original data, green is data after anomaly was introduced



# UCR\_Anomaly\_gait2\_22000\_46500\_46800.txt

Selected  
input: record  
gaitndd/hunt14 ,  
from 0:00.000 to  
5:00.000

Gait in  
Neurodegenerativ  
e Disease  
Database  
(gaitndd)

This dataset comes from someone walking on a force plate in a biomechanics lab.

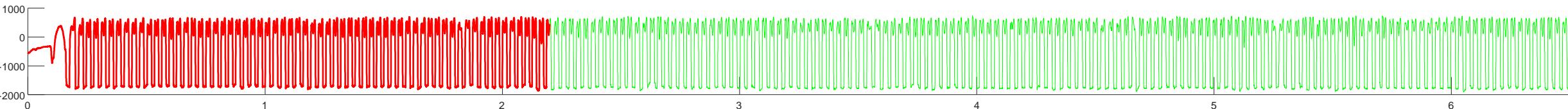
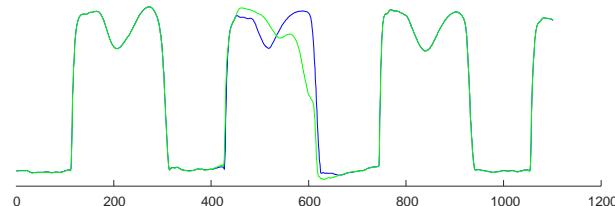
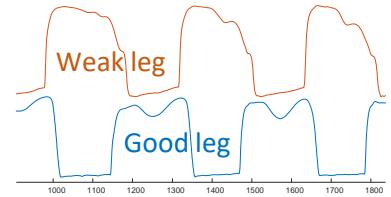
The individual had Huntington's disease, and a highly asymmetric gait (a limp).

The dataset is from the good leg (his right)

The anomaly is that at a random location we swapped in the weak leg (shifting it by a half cycle, so it lines up correctly) for one cycle.

This models an individual having a sudden, but thankfully short lived, weakness or pain in their leg.

(This is very similar to UCR\_Anomaly\_gait1\_20000\_38500\_38800.txt)



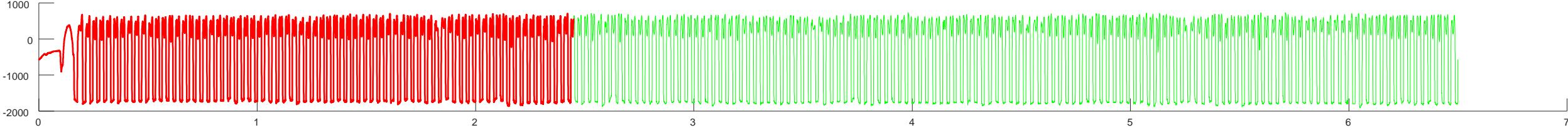
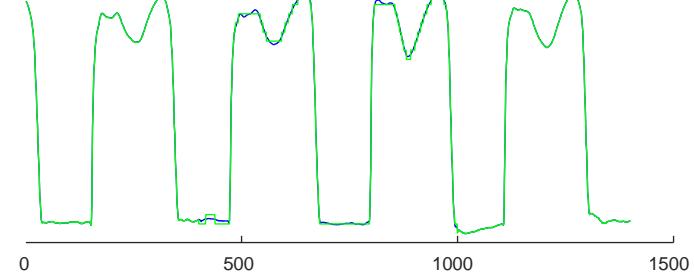
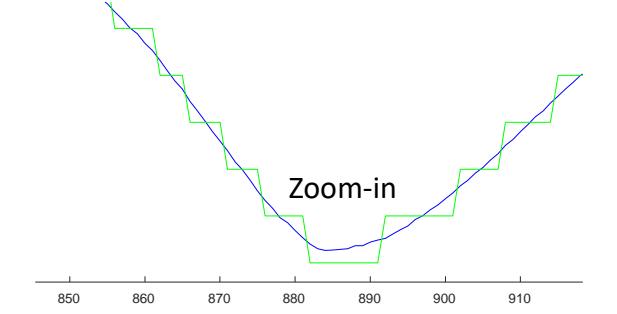
# UCR\_Anomaly\_gait3\_24500\_59900\_60500.txt

Selected  
input: record  
gaitndd/hunt14 ,  
from 0:00.000 to  
5:00.000

Gait in  
[Neurodegenerativ](#)  
[e Disease](#)  
[Database](#)  
(gaitndd)

This dataset comes from someone walking on a force plate in a biomechanics lab.  
The individual had Huntington's disease, and a highly asymmetric gait (a limp).  
The dataset is from the good leg (his right)  
The anomaly is that for two steps, we reduce the bit dept significantly

- Divide by 100
- Round
- Multiple by 100



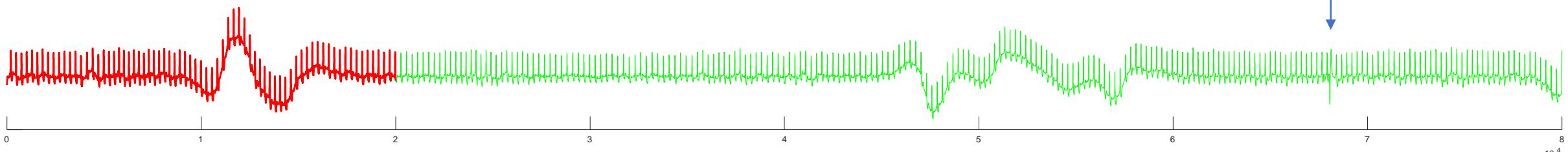
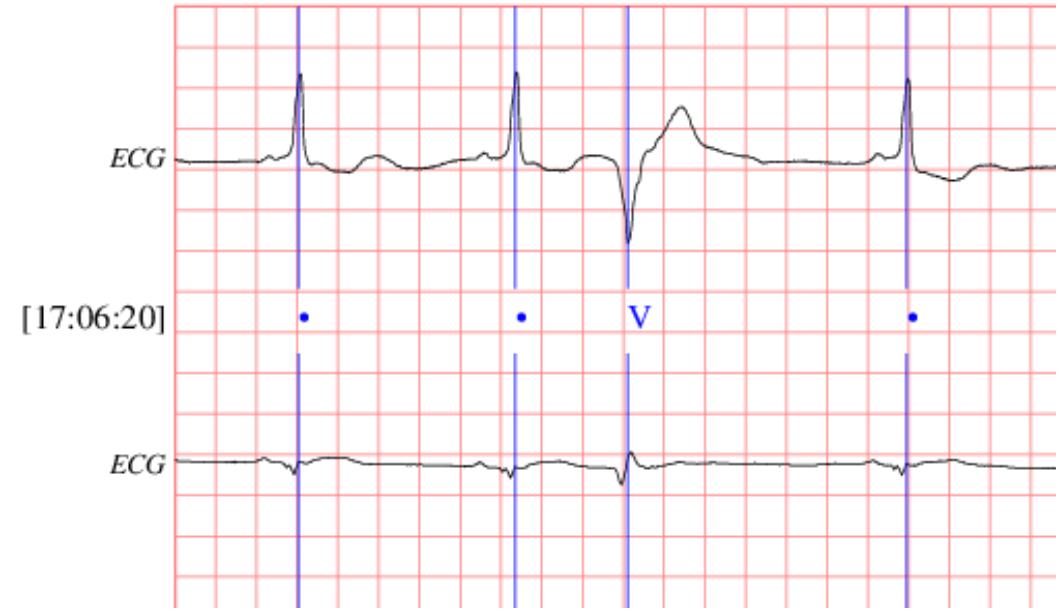
# UCR\_Anomaly\_sddb49\_20000\_67950\_68200.txt

Selected  
input: record  
sddb/49 ,  
annotator atr ,  
from  
[17:01:50.000]  
to  
[18:01:50.000]

[Sudden Cardiac  
Death Holter  
Database \(sddb\)](#)

This is a completely natural dataset. We found a long stretch of ECG that had a single anomaly (V beat).

The data was *wandering baseline*, which I suspect will confuse many algorithms.

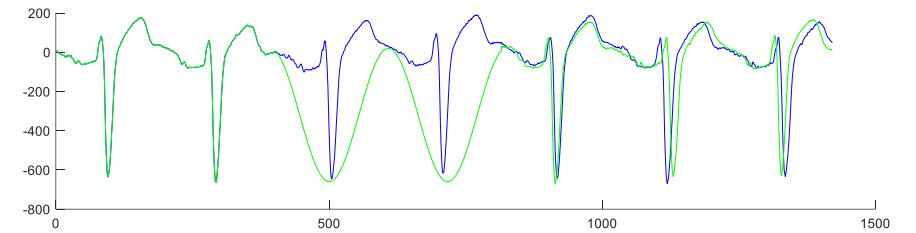
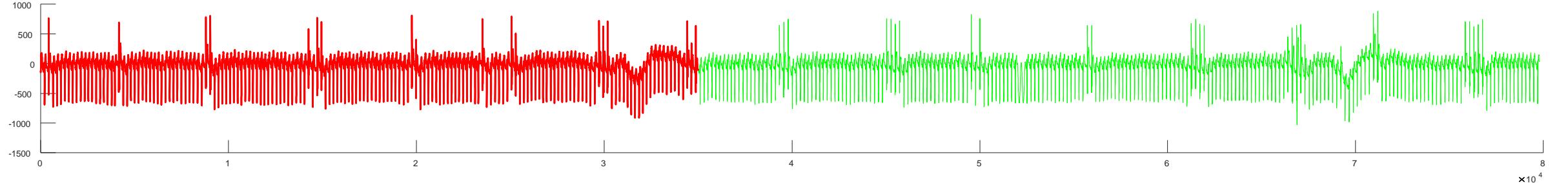


# UCR\_Anomaly\_1sddb40\_35000\_52000\_52620.txt

Selected  
input: record sddb/40 ,  
annotator atr , from [03:27:00.00 0 1] to [04:27:00.00 0 1]  
Sudden Cardiac Death Holter Database (sddb)

In this dataset there are a mixture of normal beats, and PVC beats (often in mini clusters) and wandering baselines.

We inserted 1.5 periods of a pure sine wave, with about the same mean, standard deviation and period as the real beats.

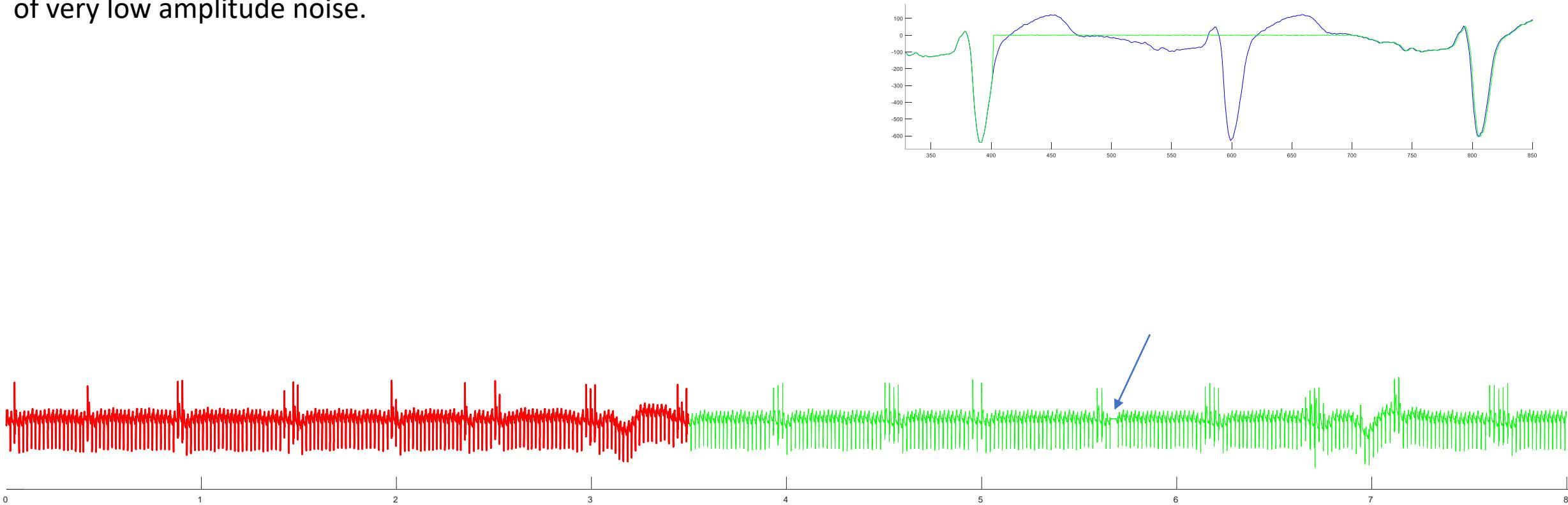


# UCR\_Anomaly\_2sddb40\_35000\_56600\_56900.txt

Selected  
input: record sddb/40 ,  
annotator atr , from  
[03:27:00.00  
0 1] to  
[04:27:00.00  
0 1]  
  
Sudden  
Cardiac  
Death  
Holter  
Database  
(sddb)  
0 1]

In this dataset there are a mixture of normal beats, and PVC beats (often in mini clusters) and wandering baselines.

We replaced 300 datapoints (about one beat) with a “constant” section of very low amplitude noise.

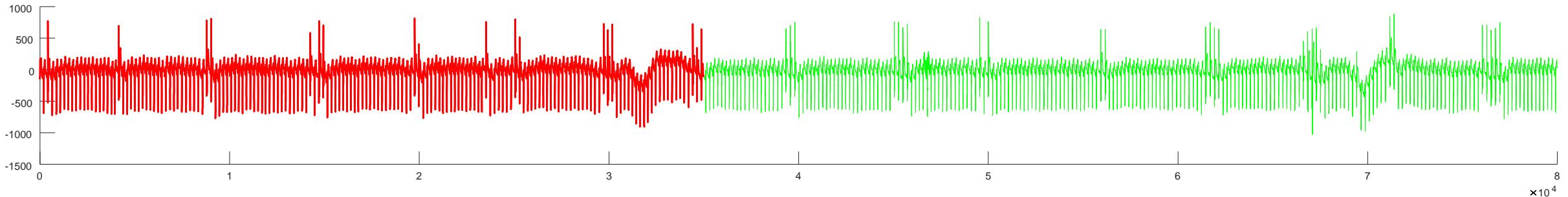
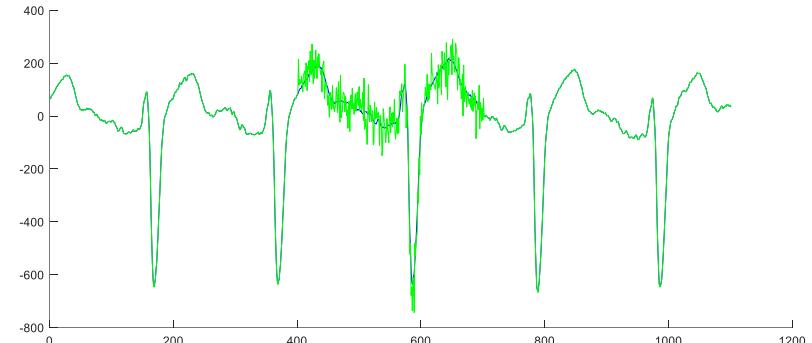


# UCR\_Anomaly\_3sddb40\_35000\_46600\_46900.txt.txt

Selected  
input: record sddb/40 ,  
annotator atr , from [03:27:00.00 0 1] to [04:27:00.00 0 1]  
Sudden Cardiac Death Holter Database (sddb)

In this dataset there are a mixture of normal beats, and PVC beats (often in mini clusters) and wandering baselines.

We added noise to a region of length 300

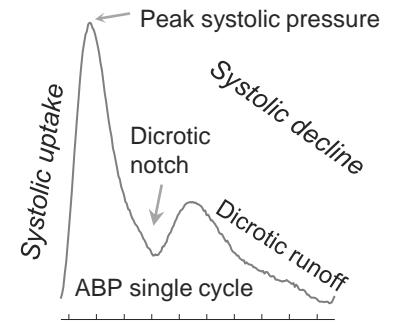
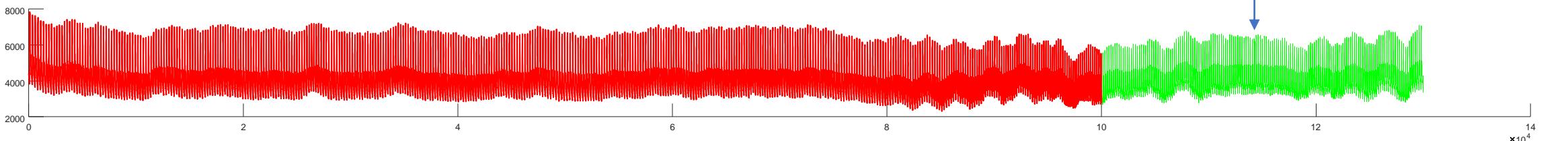


# UCR\_Anomaly\_tiltAPB1\_100000\_114283\_114350.txt

The data comes from a healthy male on a tilt table.

At first, he is supine, at around 80000, the table is tilted forward. The trace is his APB.

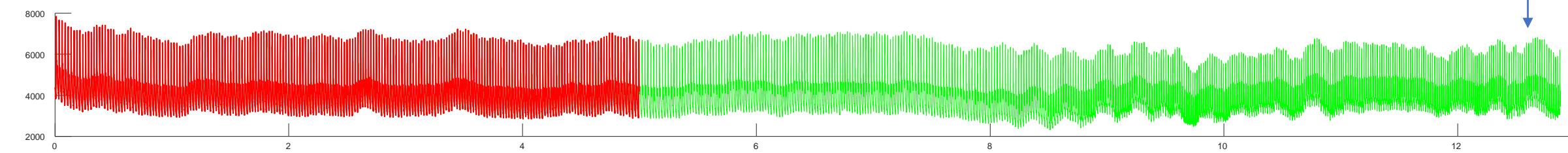
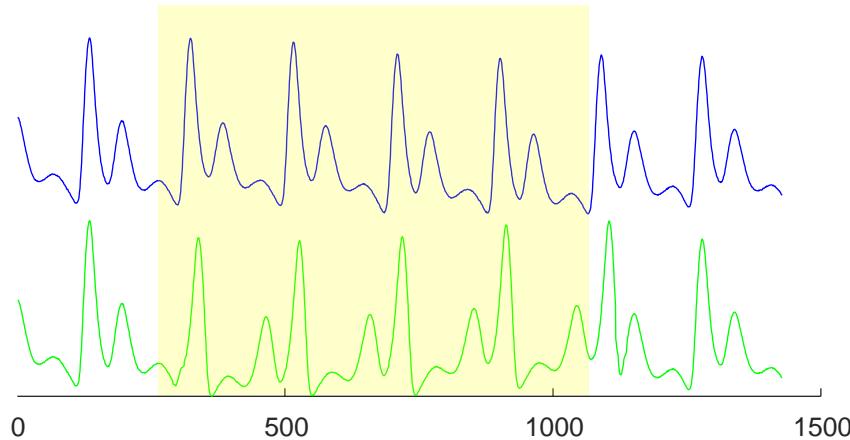
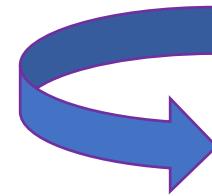
The anomaly is synthetic. There is a secondary peak after the dicrotic notch. It is normally about half the size of the peak systolic pressure. For one randomly chosen beat, we made it much greater, almost as big as the main peak.



# UCR\_Anomaly\_tiltAPB2\_50000\_124159\_124985.txt

The data comes from a healthy male on a tilt table.  
At first, he is supine, at around 80000, the table is tilted  
forward. The trace is his APB.

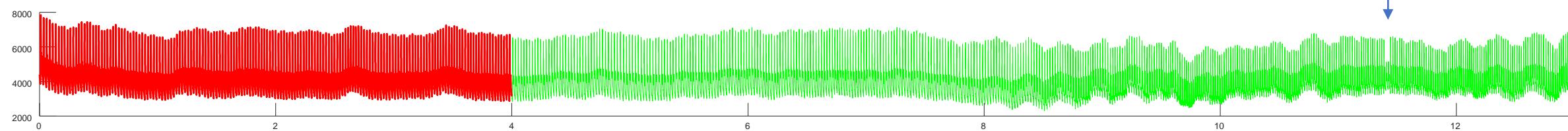
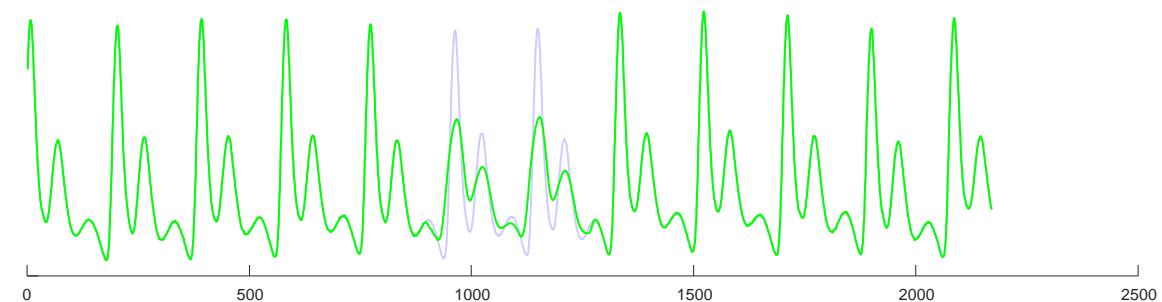
The anomaly is synthetic. We simply reverse the direction of  
the time series for about four beats.



# UCR\_Anomaly\_tiltAPB3\_40000\_114000\_114370.txt

The data comes from a healthy male on a tilt table.  
At first, he is supine, at around 80000, the table is tilted  
forward. The trace is his APB.

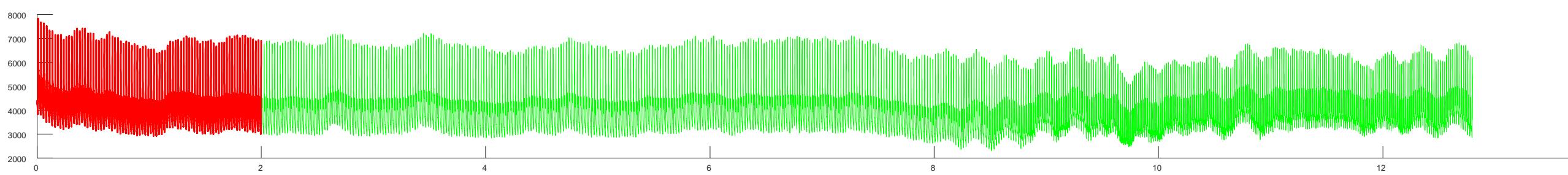
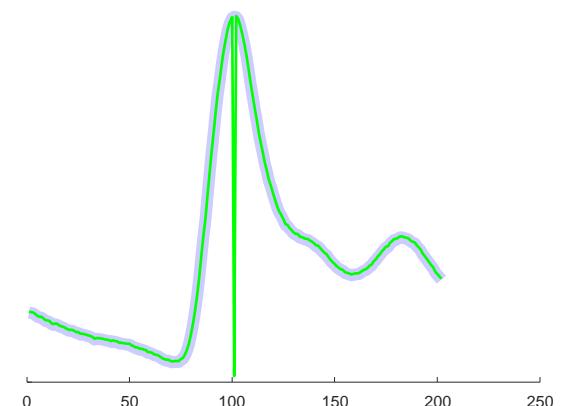
The anomaly is synthetic. We ran MATLAB's default  
smoothing on two beats. This has the effect of “dulling” the  
peaks and valleys.



# UCR\_Anomaly\_tiltAPB4\_20000\_67995\_67996.txt

The data comes from a healthy male on a tilt table.  
At first, he is supine, at around 80000, the table is tilted  
forward. The trace is his APB.

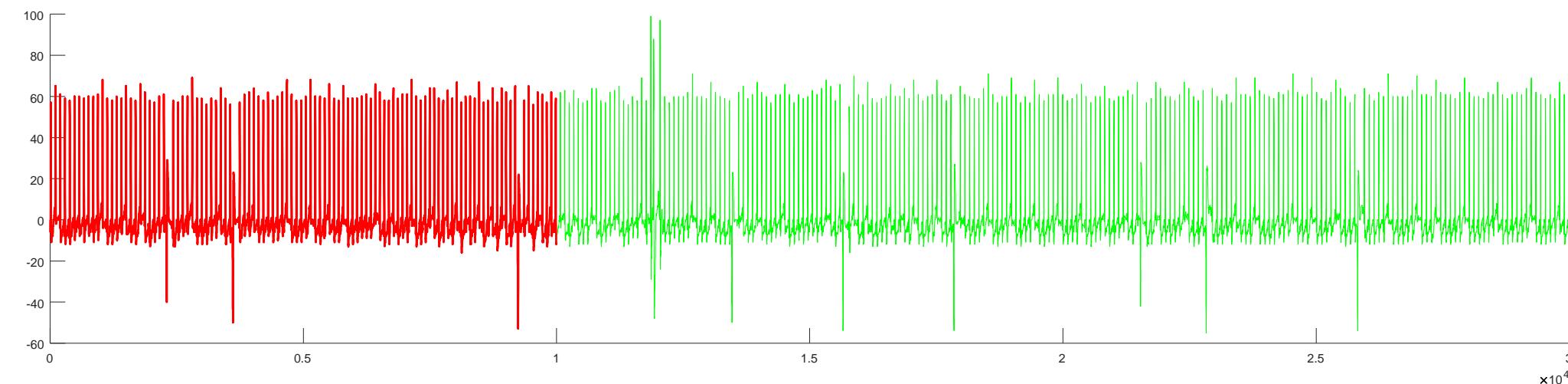
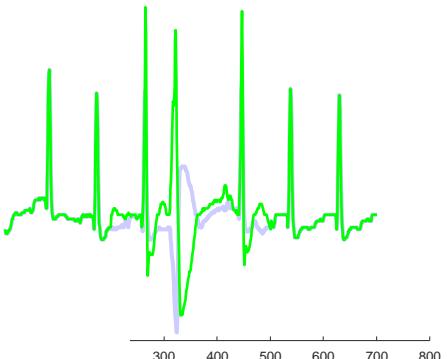
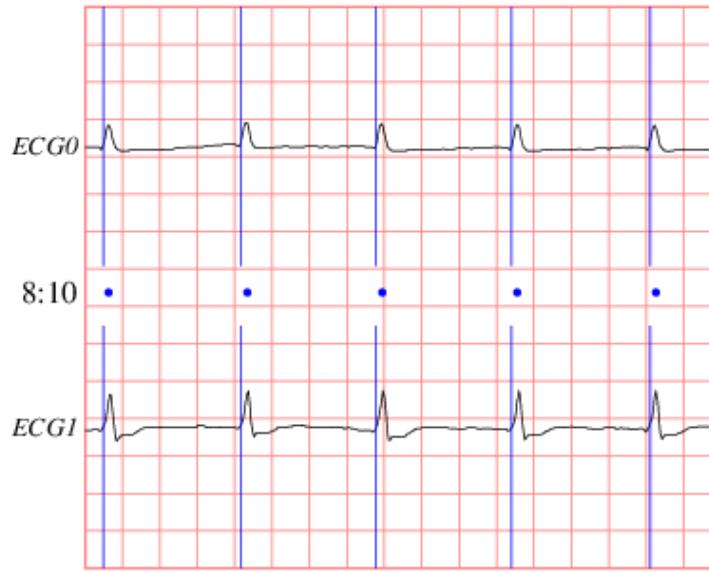
The anomaly is synthetic. We added a dropout at the top of a  
peak.



# UCR\_Anomaly\_ECG1\_10000\_11800\_12100.txt

This dataset was prepared from a two-lead ECG trace that contains some PVCs.  
Note that both train and test have multiple PVCs, that is not the anomaly.

The anomaly is synthetic. At a random location, we swapped in ECG1 for ECG2.  
That random location happened to include a PVC.



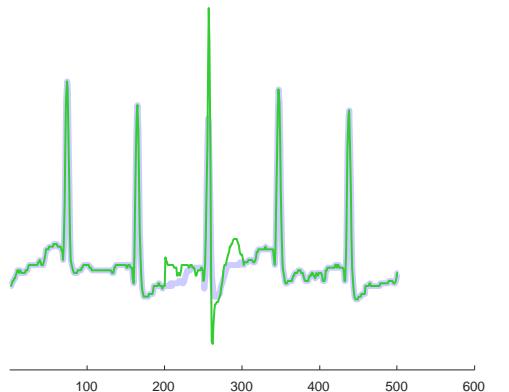
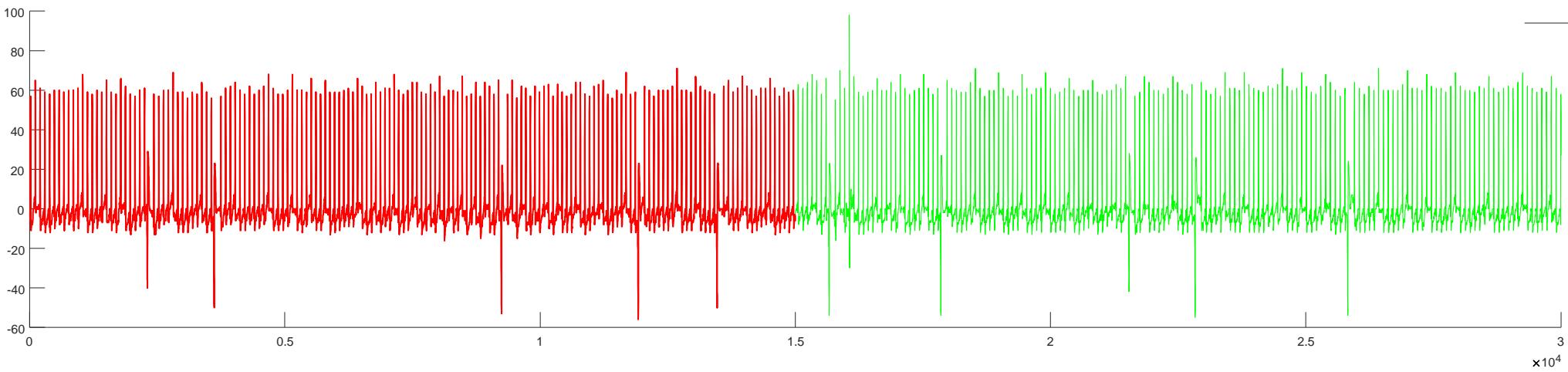
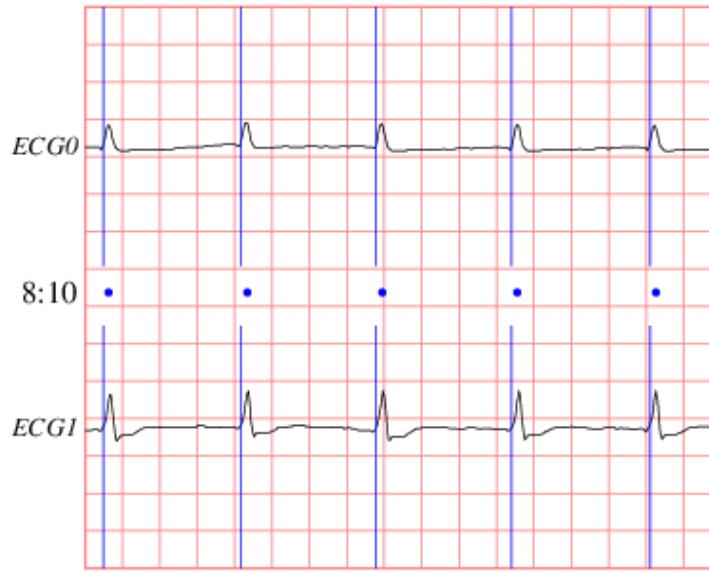
Selected  
input: record  
afpdb/n04 ,  
annotator qrs ,  
from 8:10.000 to  
8:20.000

PAF Prediction  
Challenge  
Database (afpdb)

# UCR\_Anomaly\_ECG2\_15000\_16000\_16100.txt

This dataset was prepared from a two-lead ECG trace that contains some PVCs.  
Note that both train and test have multiple PVCs, that is not the anomaly.

The anomaly is synthetic. At a random location, we swapped in ECG1 for ECG2.  
This is a more subtle version of UCR\_Anomaly\_ECG1\_10000\_11800\_12100.txt



Selected  
input: record  
afpdb/n04 ,  
annotator qrs ,  
from 8:10.000 to  
8:20.000

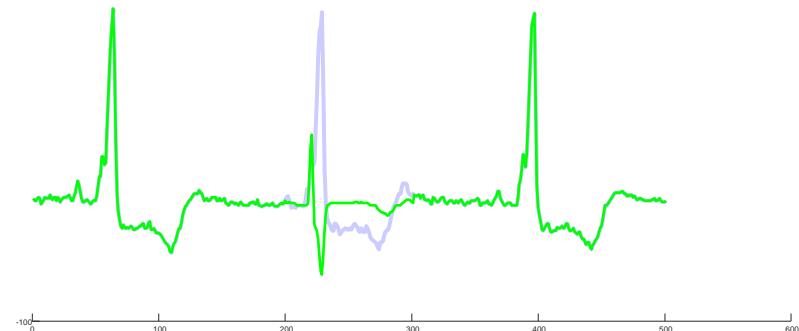
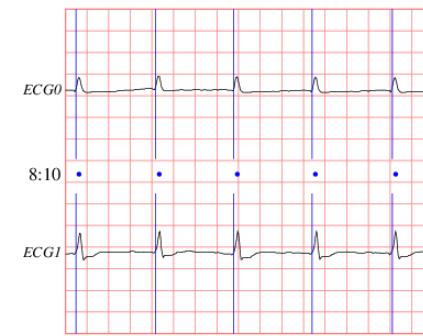
PAF Prediction  
Challenge  
Database (afpdb)

# UCR\_Anomaly\_ECG3\_8000\_17000\_17100.txt

This dataset was prepared from an ECG trace.

Note that both train and test have noise levels that wax and wane, that is not the anomaly

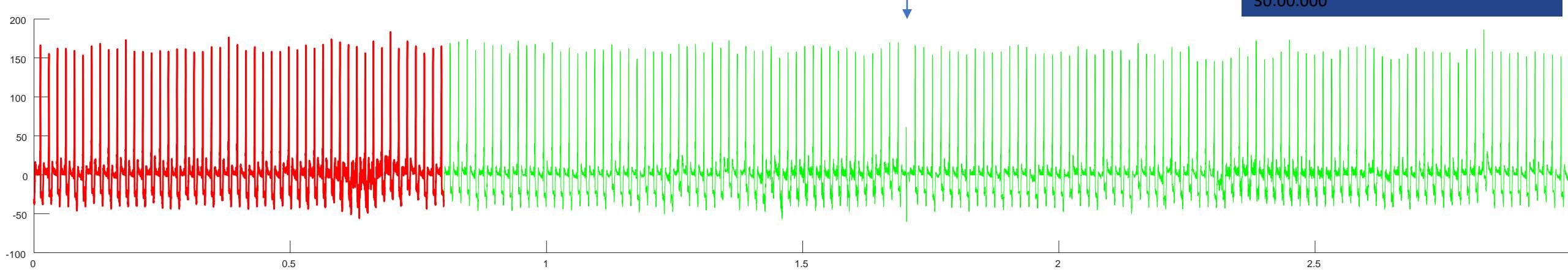
The anomaly is synthetic. At a random location, we swapped in ECG1 for ECG2.



-100 0 100 200 300 400 500 600

Selected  
input: record  
afpdb/n20 ,  
annotator qrs ,  
from 9:50.000 to  
30:00.000

[PAF Prediction  
Challenge  
Database \(afpdb\)](#)

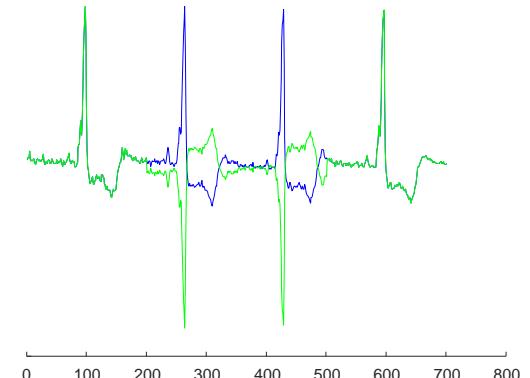


# UCR\_Anomaly\_ECG4\_5000\_16800\_17100.txt

This dataset was prepared from an ECG trace.

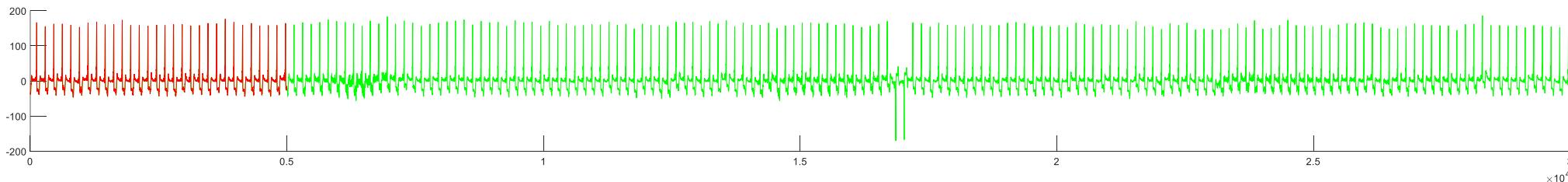
Note that both train and test have noise levels that wax and wane, that is not the anomaly

The anomaly is synthetic. At a random location, we turned the data, about two beats worth, upside down (by multiplying the subsequence by -1)



Selected  
input: record  
afpdb/n20 ,  
annotator qrs ,  
from 9:50.000 to  
30:00.000

[PAF Prediction Challenge Database \(afpdb\)](#)



# UCR\_Anomaly\_respiration1\_100000\_110260\_110412.txt

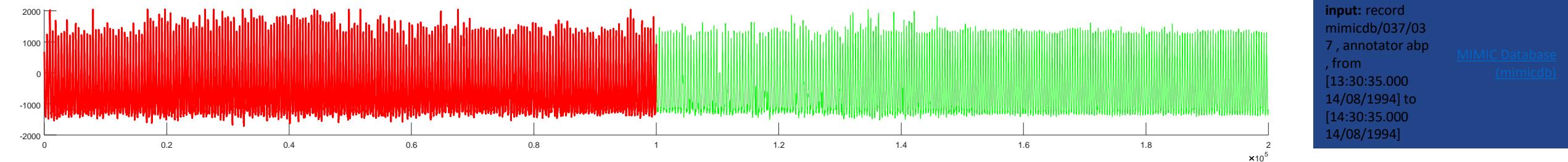
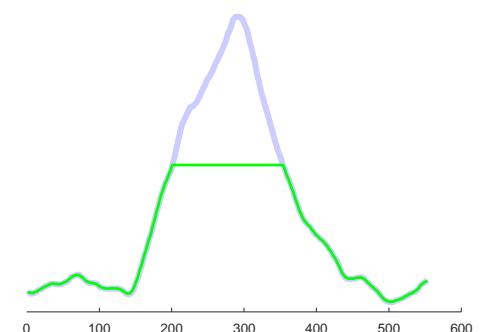
This is a respiration dataset.

Many of these datasets have hard upper and lower limits, due to

- 1) Bit depth: Here we see a few peaks maximized at 2047 (which is  $2^{12}/2 - 1$ )
- 2) Mechanical limits of the elasticity of the chest band (hardish limits, produces *almost* constant values)

The mechanical limits can change with movement, or with a nurse placing her hand under the band for some reason etc.

Here we model such an event, for a random breath cycle, the maximum value was hard limited to near zero.



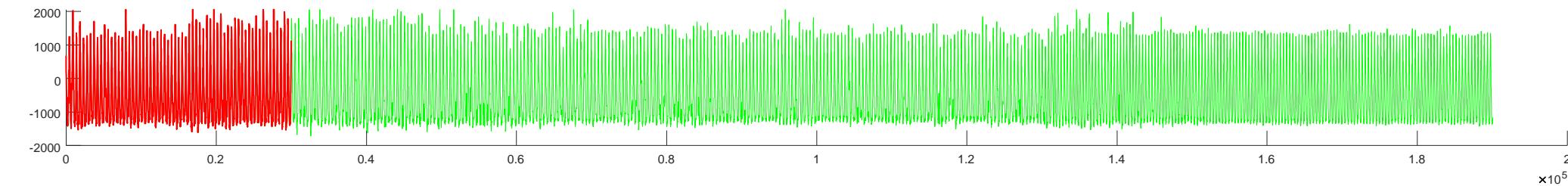
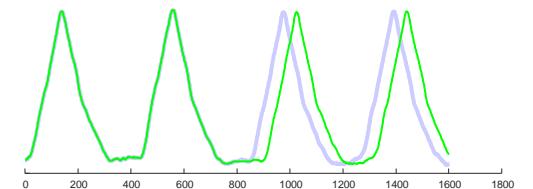
# UCR\_Anomaly\_resperation2\_30000\_168250\_168251.txt

This is a respiration dataset.

The breath cycles here are very regular (implying a deep sleep).

Here we created an anomaly by increasing the time the patient held their breath (by a subtle amount).

This models sleep apnea



**Selected**  
input: record  
mimicdb/037/03  
7 , annotator abp  
, from  
[13:30:35.000  
14/08/1994] to  
[14:30:35.000  
14/08/1994]

MIMIC Database  
(mimicdb)

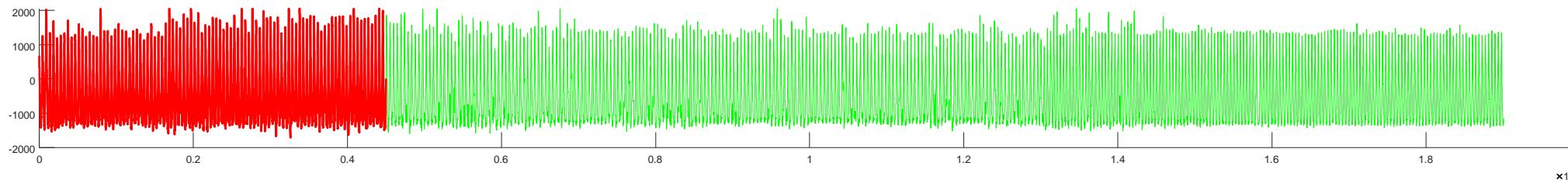
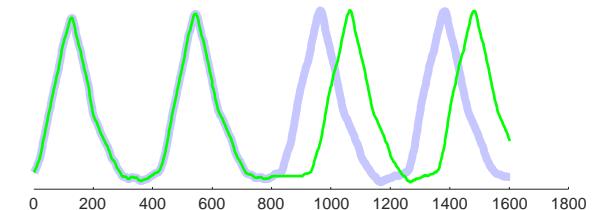
# UCR\_Anomaly\_resperation3\_45000\_158250\_158251.txt

This is a respiration dataset.

The breath cycles here are very regular (implying a deep sleep).

Here we created an anomaly by increasing the time the patient held their breath (by a significant amount).

This models sleep apnea



**Selected**  
input: record  
mimicdb/037/03  
7 , annotator abp  
, from  
[13:30:35.000  
14/08/1994] to  
[14:30:35.000  
14/08/1994]

MIMIC Database  
(mimicdb)

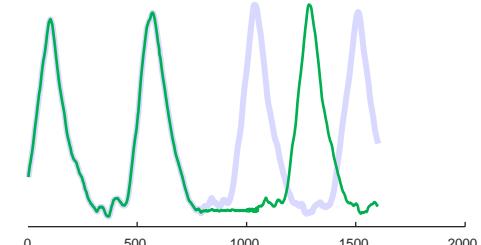
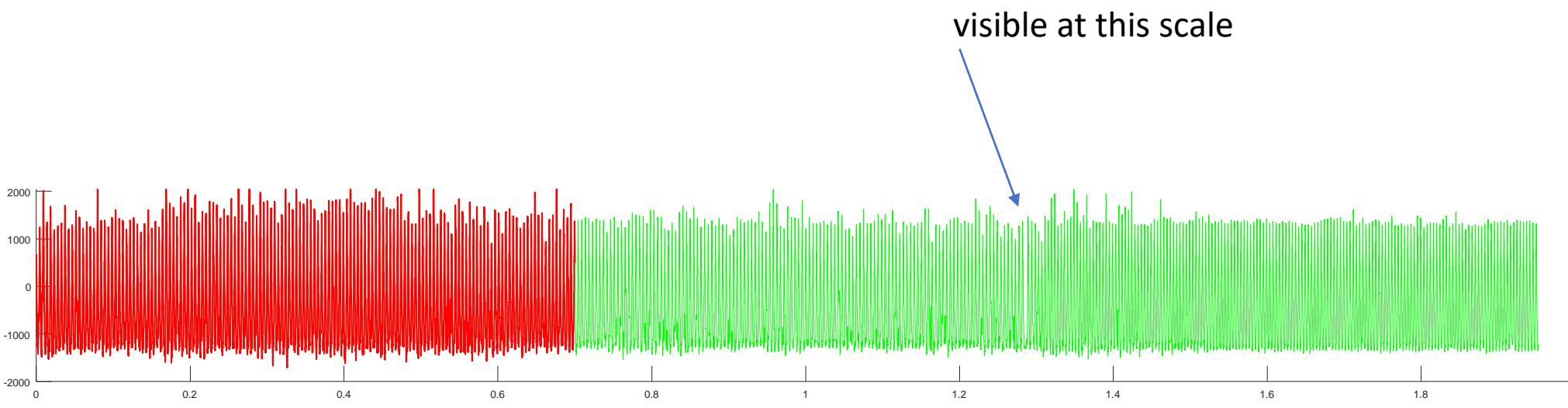
# UCR\_Anomaly\_respiration4\_70000\_128430\_128431.txt

This is a respiration dataset.

The breath cycles here are very regular (implying a deep sleep).

Here we created an anomaly by increasing the time the patient held their breath (by a very large amount).

This models sleep apnea



Selected  
input: record  
mimicdb/037/03  
7 , annotator abp  
, from  
[13:30:35.000  
14/08/1994] to  
[14:30:35.000  
14/08/1994]

MIMIC Database  
(mimicdb)

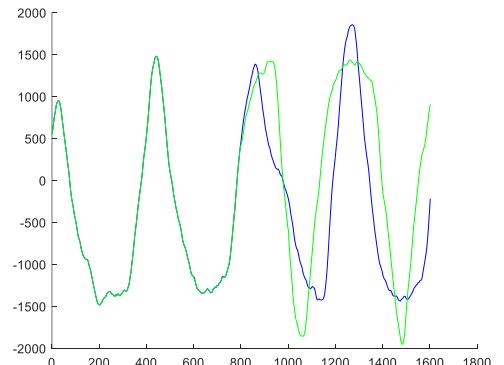
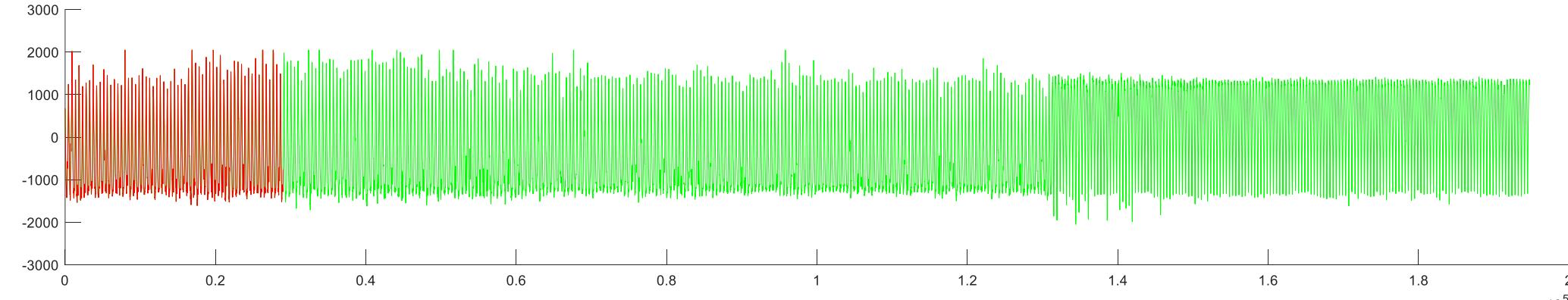
# UCR\_Anomaly\_resperation8\_29000\_131230\_131231.txt

This is a respiration dataset.

The breath cycles here are very regular (implying a deep sleep).

Here we created an anomaly by simply turning the data upside down at a random location.

We cleaned up and smoothed the “join” so as to make it less visible.



**Selected**  
input: record  
mimicdb/037/03  
7 , annotator abp  
, from  
[13:30:35.000  
14/08/1994] to  
[14:30:35.000  
14/08/1994]

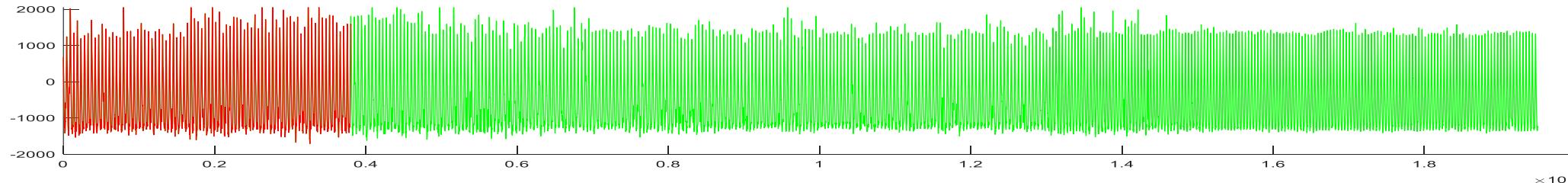
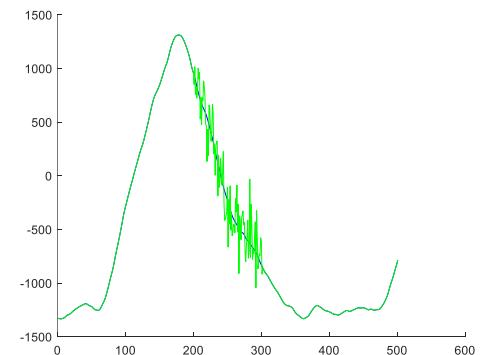
MIMIC Database  
(mimicdb)

# UCR\_Anomaly\_resperation9\_38000\_143411\_143511.txt

This is a respiration dataset.

The breath cycles here are very regular (implying a deep sleep).

Here we created an anomaly by adding noise.



**Selected**  
**input:** record  
mimicdb/037/03  
7 , annotator abp  
, from  
[13:30:35.000  
14/08/1994] to  
[14:30:35.000  
14/08/1994]

MIMIC Database  
(mimicdb)

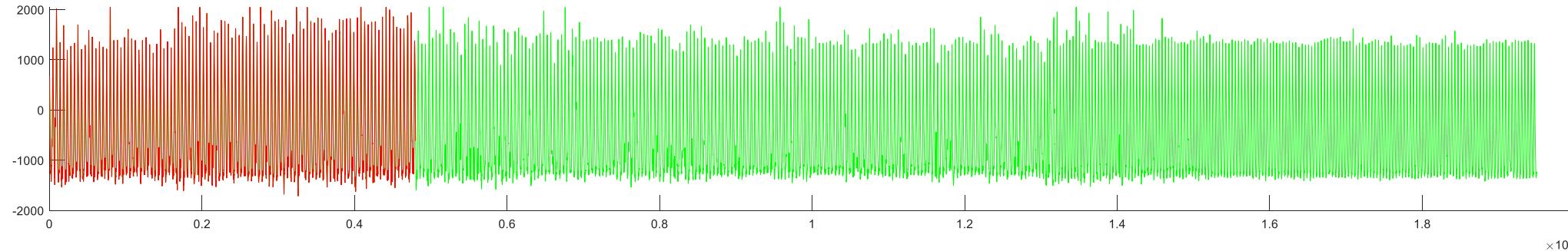
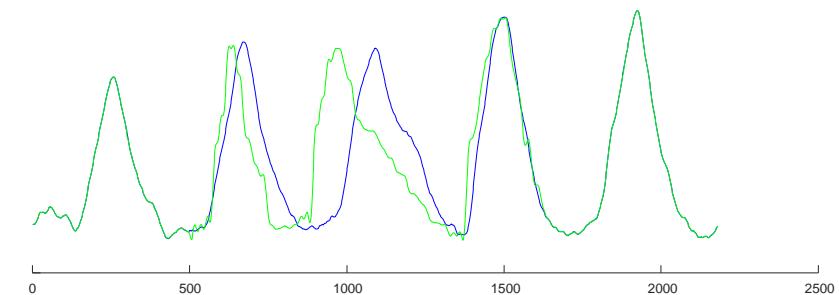
# UCR\_Anomaly\_resperation10\_48000\_130700\_131880.txt

This is a respiration dataset.

The breath cycles here are very regular (implying a deep sleep).

Here we created an anomaly by adding time warping.

```
function [warped_T] = add_warping_one_time_series(T,p)
i = randperm(length(T));
i = sort(i(1:end-floor(length(T) * p)));
warped_T = smooth(resample(T(i),length(T),length(i)),1);
```



**Selected**  
input: record  
mimicdb/037/03  
7 , annotator abp  
, from  
[13:30:35.000  
14/08/1994] to  
[14:30:35.000  
14/08/1994]

MIMIC Database  
(mimicdb)

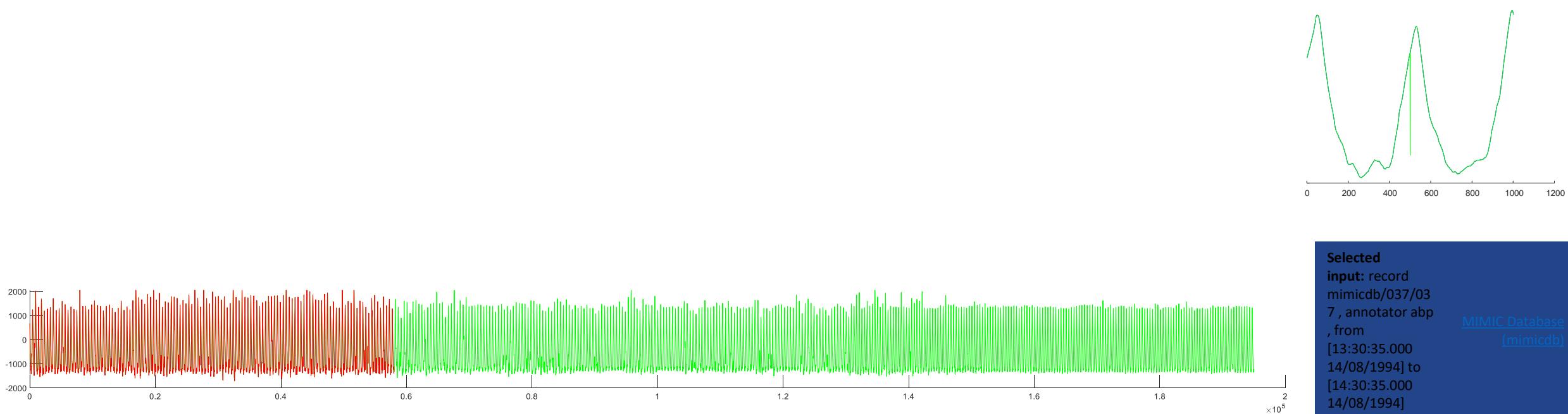
# UCR\_Anomaly\_resperation11\_58000\_110800\_110801.txt

This is a respiration dataset.

The breath cycles here are very regular (implying a deep sleep).

In some legacy systems, missing values are encoded at -999.

Here we randomly inserted such a value.

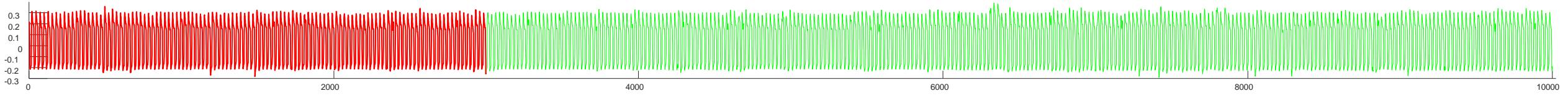
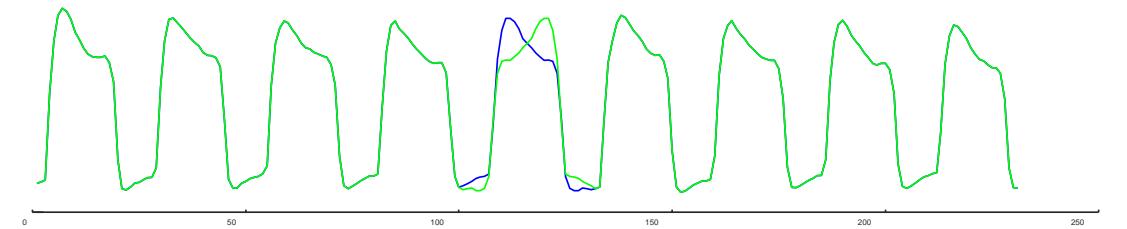


# UCR\_Anomaly\_insectEPG1\_3000\_7000\_7030.txt

This is an insect EPG dataset, an Asian Citrus Psyllid.

This is a trace of insect feeding behavior.

We reversed a single cycle

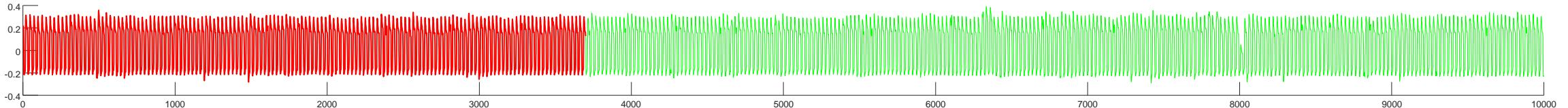
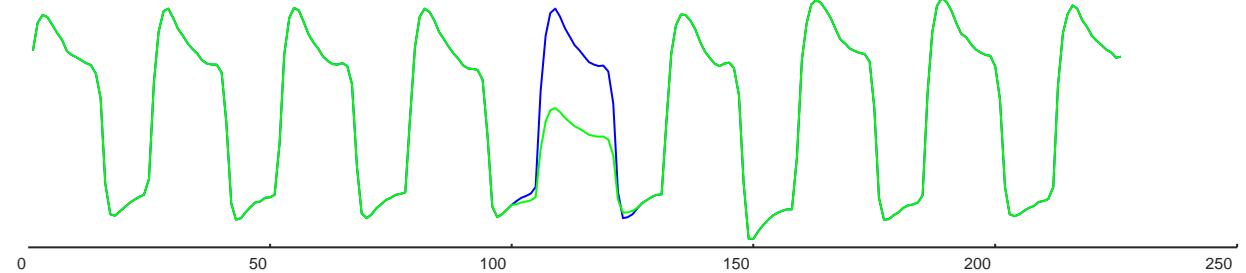


# UCR\_Anomaly\_insectEPG2\_3700\_8000\_8025.txt

This is an insect EPG dataset, an Asian Citrus Psyllid.

This is a trace of insect feeding behavior.

We reduced the amplitude of a single cycle by half

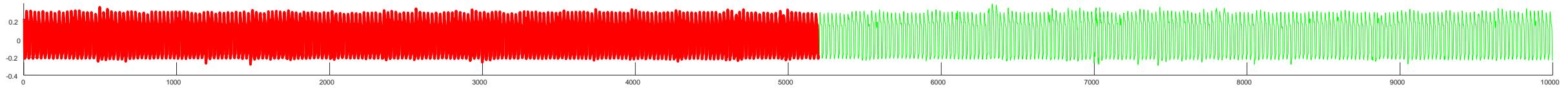
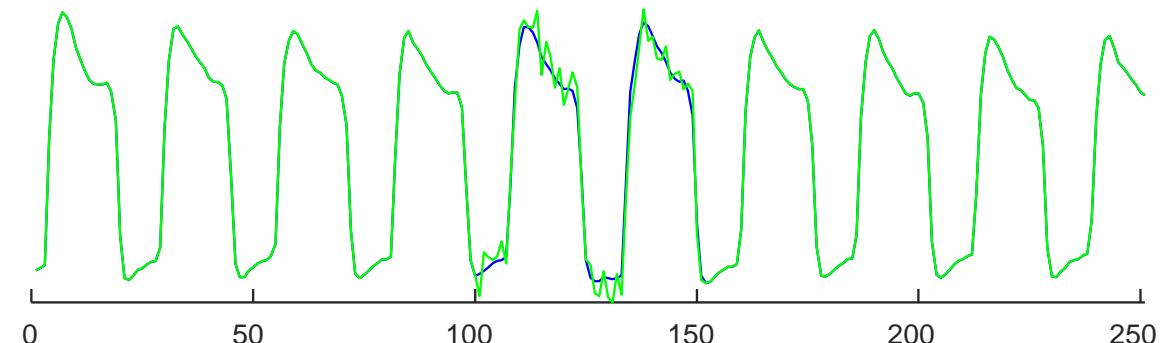


# UCR\_Anomaly\_insectEPG3\_5200\_7000\_7050.txt

This is an insect EPG dataset, an Asian Citrus Psyllid.

This is a trace of insect feeding behavior.

We added noise to two cycles

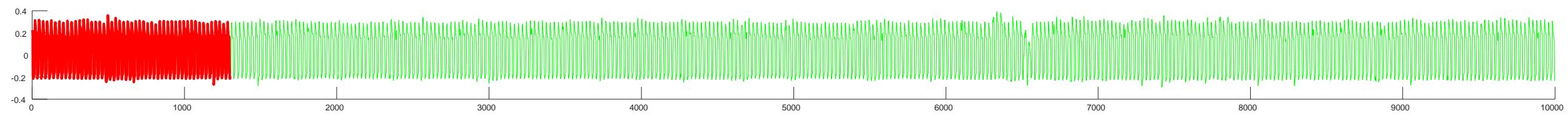
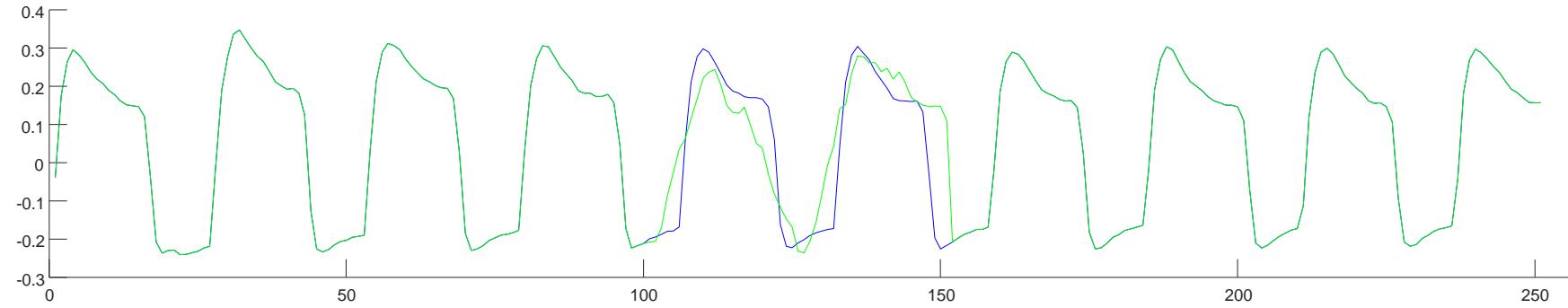


# UCR\_Anomaly\_insectEPG4\_1300\_6508\_6558.txt

This is an insect EPG dataset, an Asian Citrus Psyllid.

This is a trace of insect feeding behavior.

We added distorted two cycles by adding noise, then smoothing.

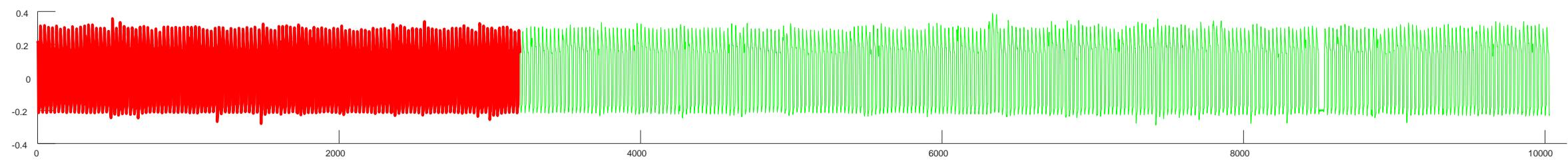
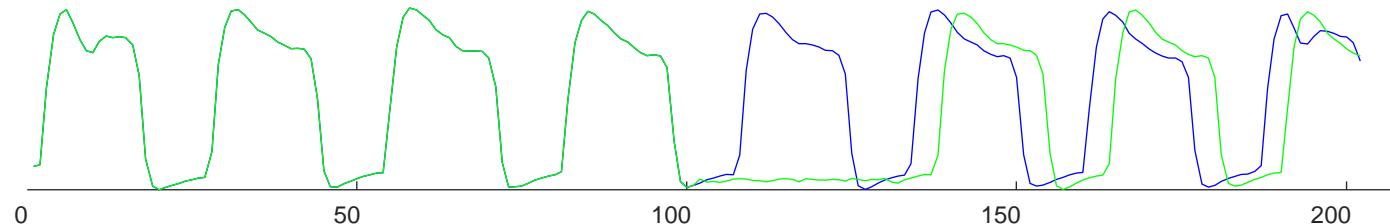


# UCR\_Anomaly\_insectEPG5\_3200\_8500\_8501.txt

This is an insect EPG dataset, an Asian Citrus Psyllid.

This is a trace of insect feeding behavior.

We added some space between two cycles



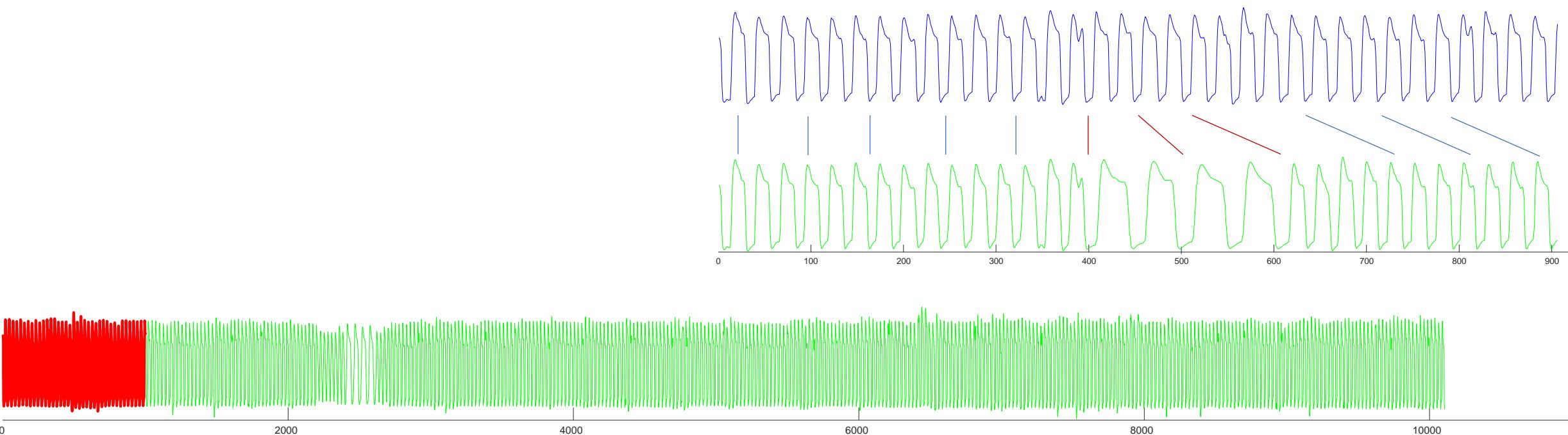
# UCR\_Anomaly\_insectEPG6\_1000\_2400\_2505.txt

This is an insect EPG dataset, an Asian Citrus Psyllid.

This is a trace of insect feeding behavior.

For four cycles, we slowed down the behavior by a factor of two, slightly smoothing the result to remove interpolation artifacts.

```
T = [ T(1:start_anomaly); T(start_anomaly:0.5:end_anomaly) ; T(end_anomaly:end) ];
T(start_anomaly-200:end_anomaly+200) = smooth(T(start_anomaly-200:end_anomaly+200),1);
```



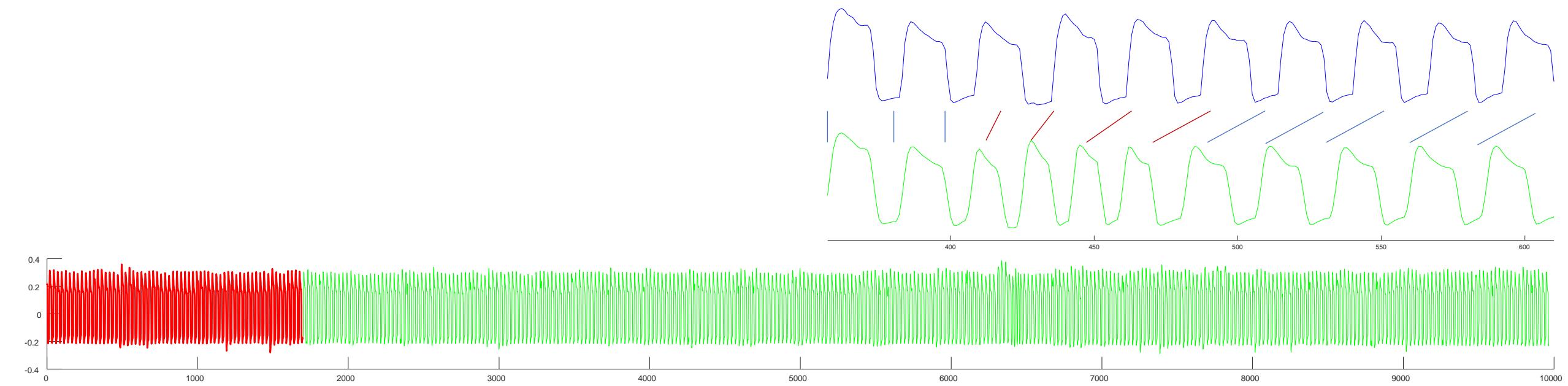
# UCR\_Anomaly\_insectEPG7\_1700\_6400\_6505.txt

This is an insect EPG dataset, an Asian Citrus Psyllid.

This is a trace of insect feeding behavior.

For four cycles, we speed up the behavior by a factor of 50%, slightly smoothing the result to remove interpolation artifacts.

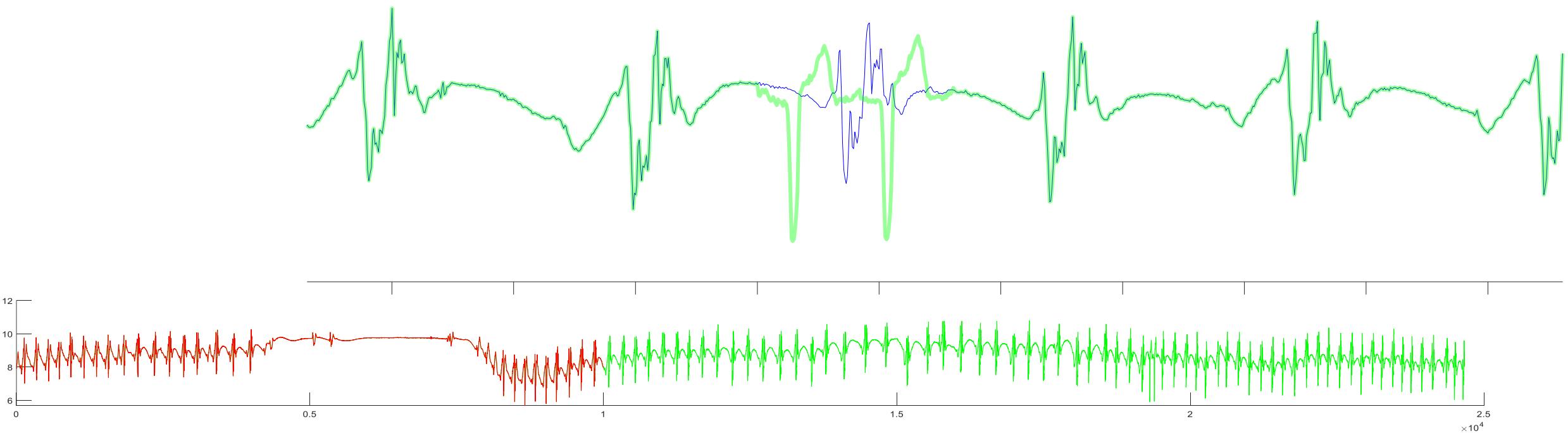
```
T = [ T(1:start_anomaly);    T(start_anomaly:1.5:end_anomaly) ;    T(end_anomaly:end)    ];
T(start_anomaly-200:end_anomaly+200) = smooth(T(start_anomaly-200:end_anomaly+200),1);
```



# UCR\_Anomaly\_MesoplodonDensirostris\_10000\_19280\_19440.txt

This data comes from an accelerometer attached to a Blainville's beaked whale.

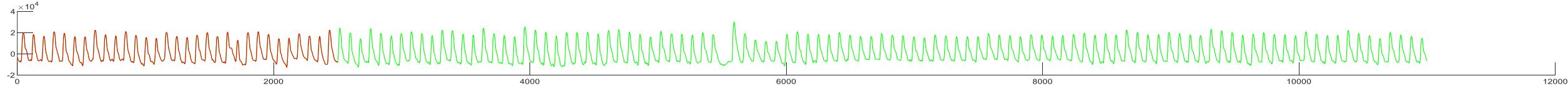
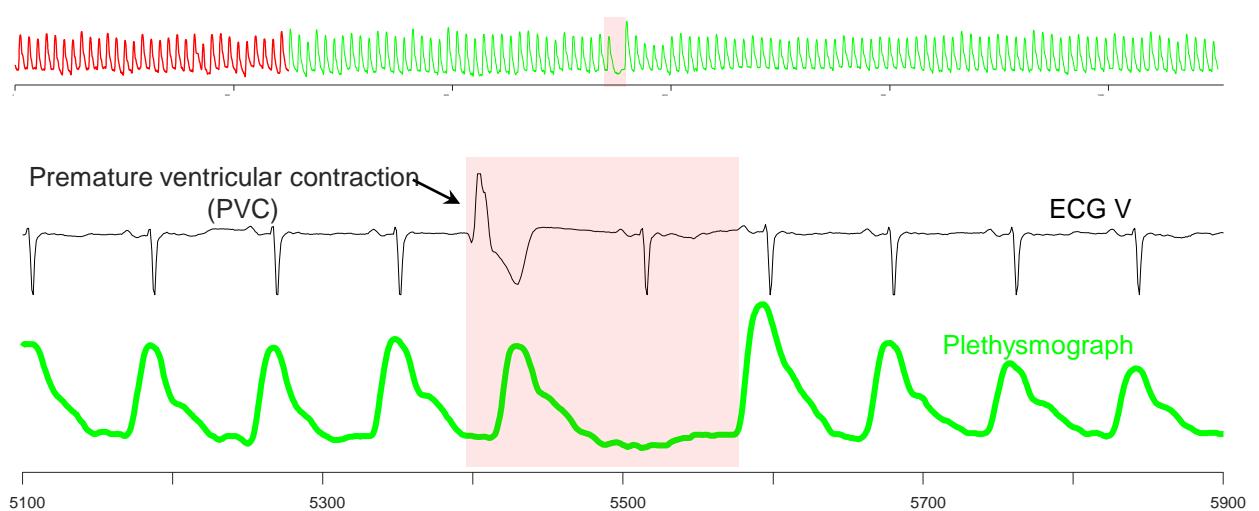
We replace one swim cycle with two human heartbeats, scaled (*very approximately*) to have about the same mean and variance as the data they replaced



# UCR\_Anomaly\_BIDMC1\_2500\_5400\_5600.txt

**Selected input:** record  
bidmc/bidmc01 , annotator ann1 ,  
from 0:00.000 to 1:00.000

Here the anomaly is a little subtle, how can we be so confident that is it semantically an anomaly? We can make this assertion because we examined the electrocardiogram that was recorded in parallel. This was the only region that had an abnormal heartbeat, a PVC. Note that there is a slight lag in the timing, as an ECG is an electrical signal, and the pleth signal is mechanical. However, the scoring functions typically have a little “play” to avoid the brittleness of requiring spurious precision.



# UCR\_Anomaly\_park3m\_60000\_72150\_72495.txt

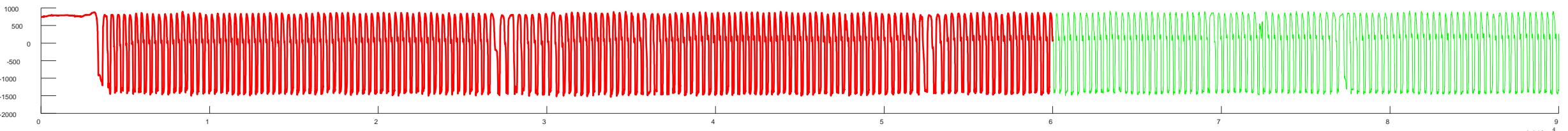
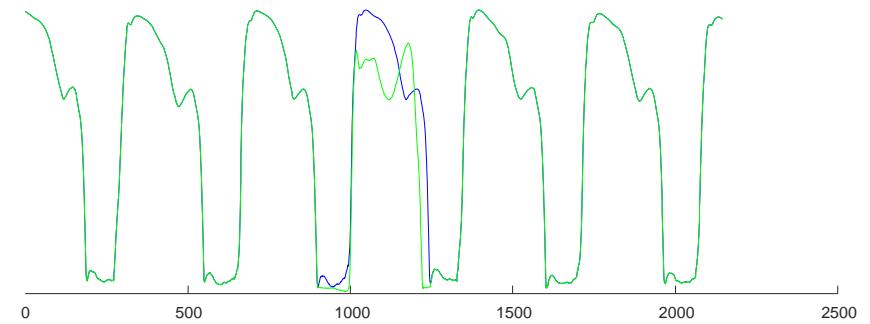
We took the right foot of a person with an asymmetric gait because of Parkinson's disease, and swapped in one cycle from his left foot,



**Selected**  
**input:** record  
gaitndd/park3 ,  
from 0:00.000 to  
5:00.000

Gait in

[Neurodegenerati  
ve Disease  
Database  
\(gaitndd\)](#)



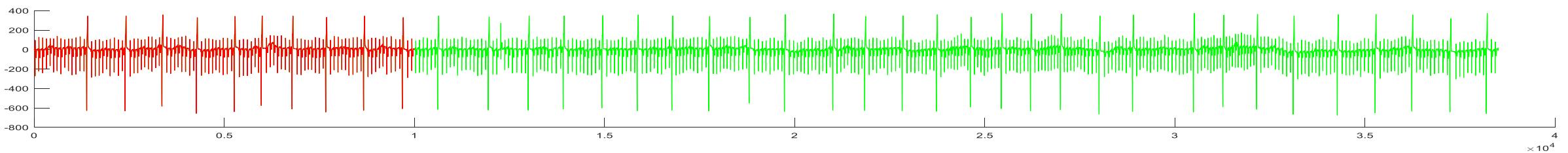
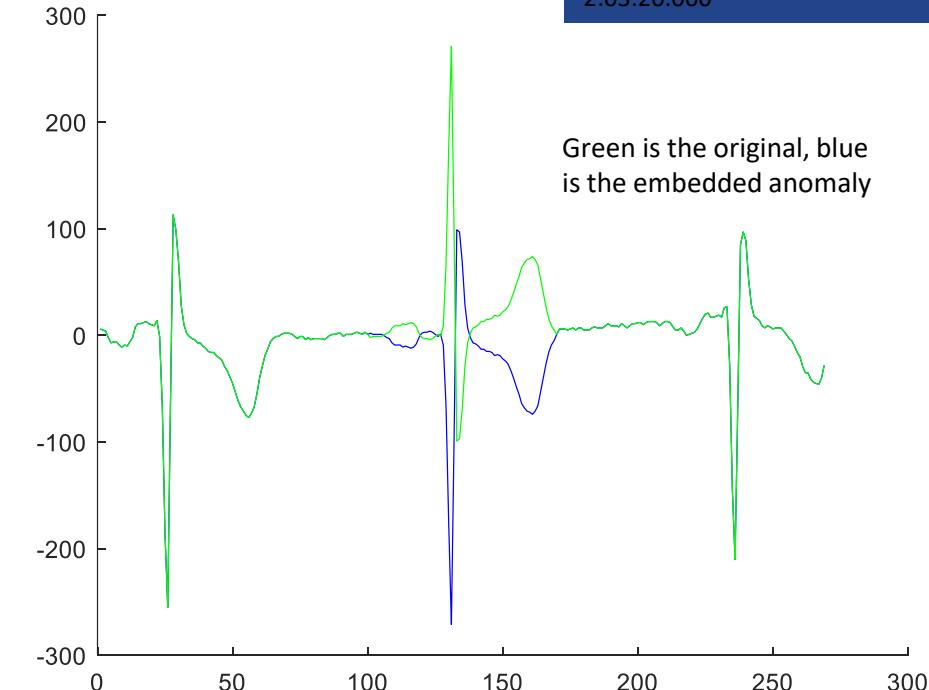
# UCR\_Anomaly\_apneaecg\_10000\_12240\_12308.txt

Selected  
input: record  
apnea-  
ecg/b05 ,  
annotator  
apn , from  
1:03:20.000  
to  
2:03:20.000

[Apnea-ECG  
Database  
\(apnea-ecg\)](#)

This dataset contains a mixture of normal beats, with approximately every tenth beat being a Premature ventricular contraction.

Here we simply flipped one normal beat upside down.



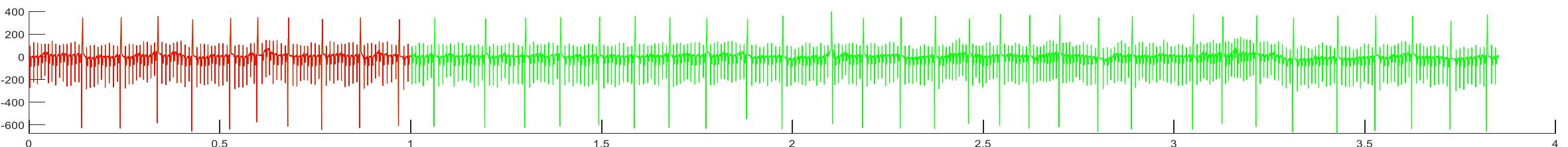
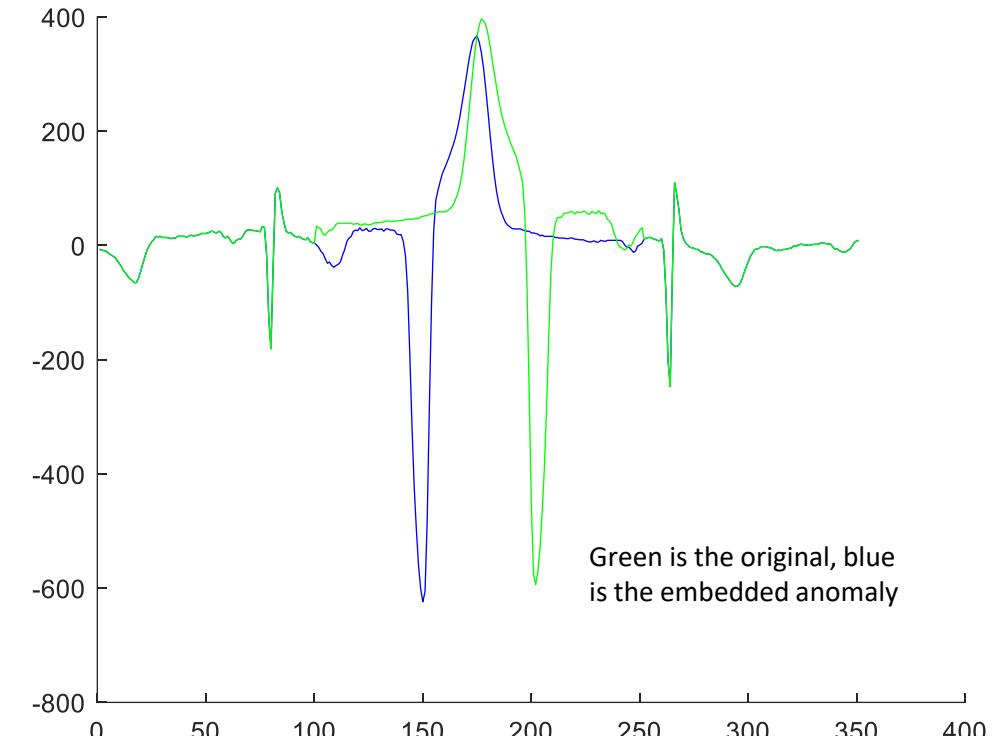
# UCR\_Anomaly\_apneaecg2\_10000\_20950\_21100.txt

Selected  
input: record  
apnea-  
ecg/b05 ,  
annotator  
apn , from  
1:03:20.000  
to  
2:03:20.000

[Apnea-ECG  
Database  
\(apnea-ecg\)](#)

This dataset contains a mixture of normal beats, with approximately every tenth beat being a Premature ventricular contraction.

Here we simply flipped one PVC beat left to right down.



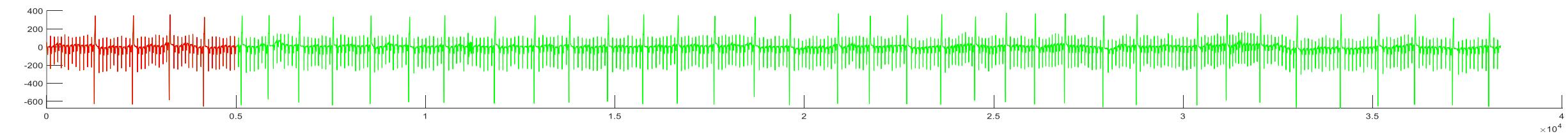
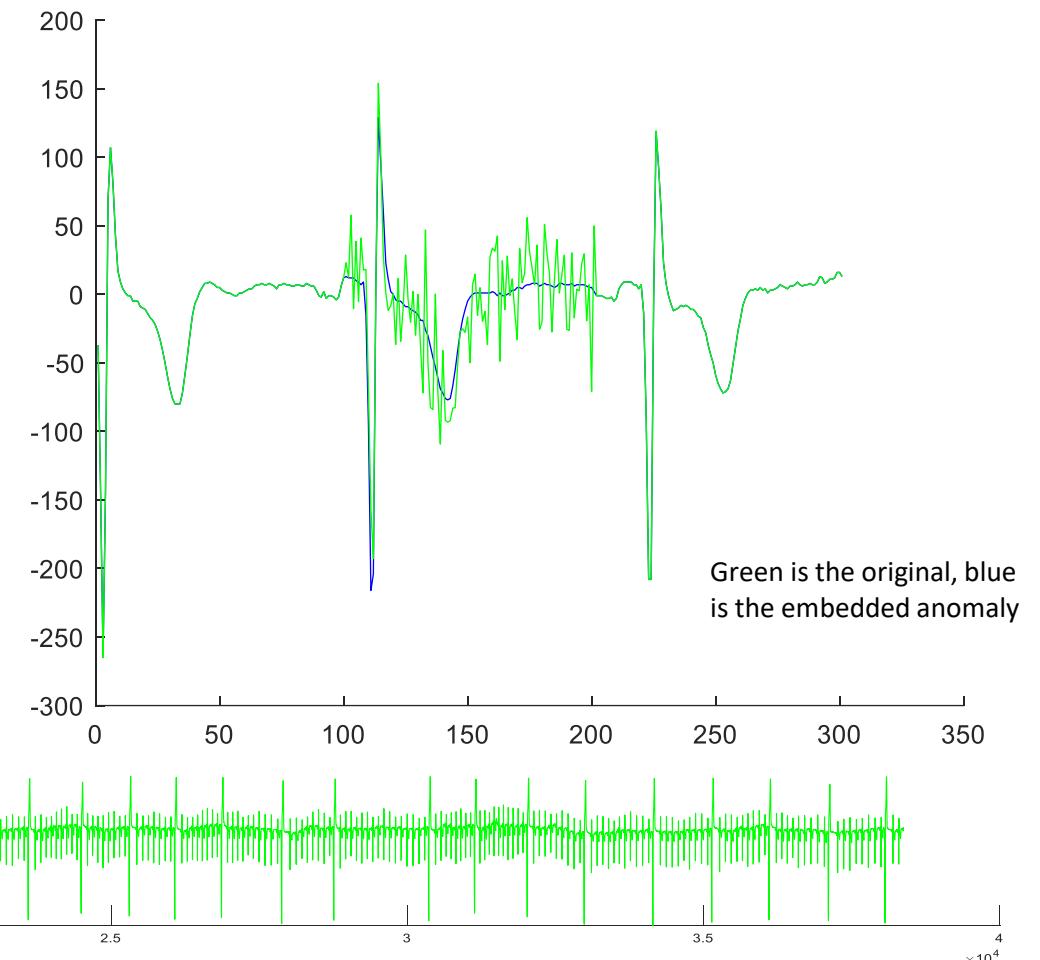
# UCR\_Anomaly\_apneaecg3\_5000\_11111\_11211.txt

Selected  
input: record  
apnea-  
ecg/b05 ,  
annotator  
apn , from  
1:03:20.000  
to  
2:03:20.000

[Apnea-ECG  
Database  
\(apnea-ecg\)](#)

This dataset contains a mixture of normal beats, with *approximately* every tenth beat being a Premature ventricular contraction.

Here we simply added noise to a short region.



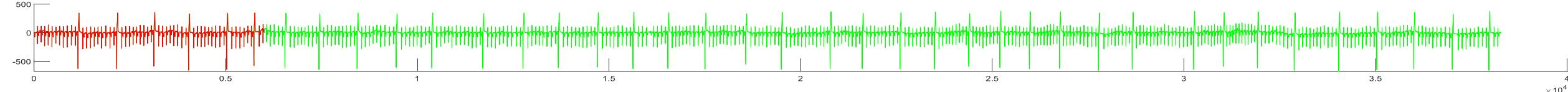
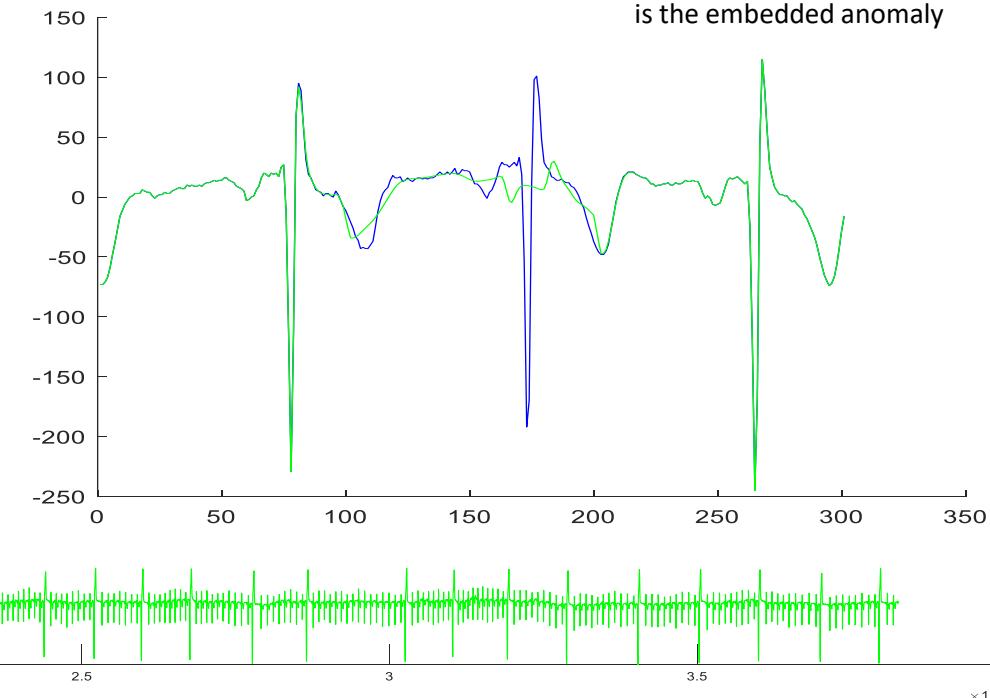
# UCR\_Anomaly\_apneaecg4\_6000\_16000\_16100.txt

Selected  
input: record  
apnea-  
ecg/b05 ,  
annotator  
apn , from  
1:03:20.000  
to  
2:03:20.000

[Apnea-ECG  
Database  
\(apnea-ecg\)](#)

This dataset contains a mixture of normal beats, with *approximately* every tenth beat being a Premature ventricular contraction.

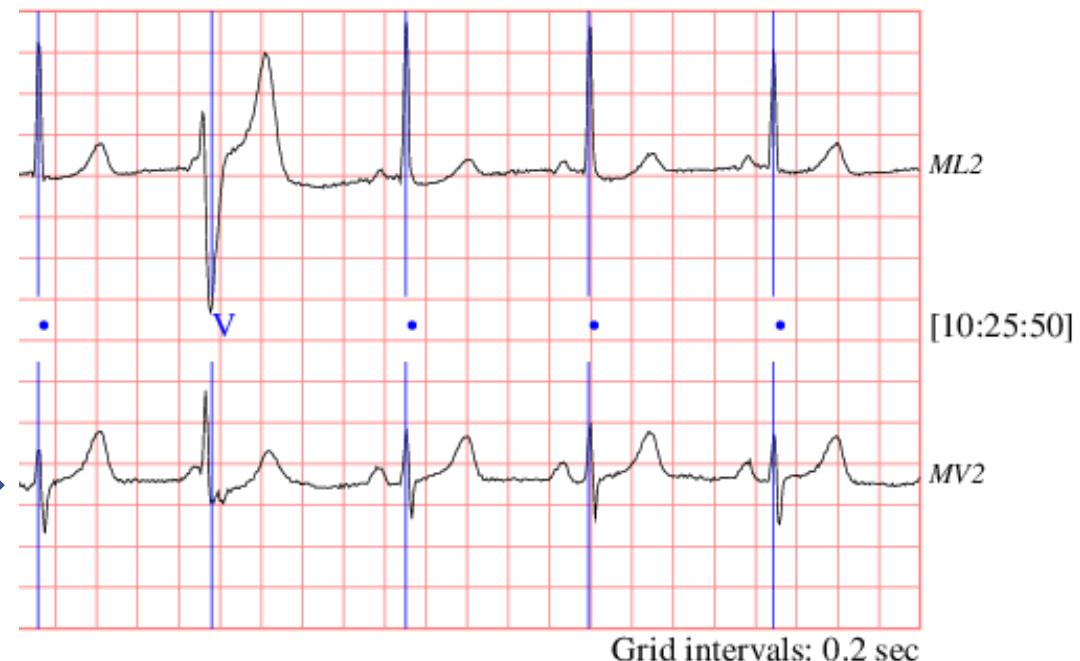
Here we *oversmoothed* a single beat. This tends to depress the peak and valley



# UCR\_Anomaly\_s20101m\_10000\_35774\_35874.txt

This dataset contains normal beats, with a single ventricular beat between 35774 and 35874.

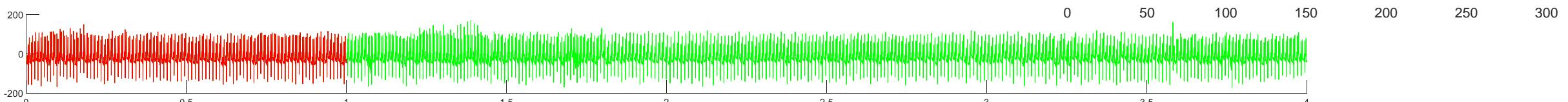
Here we use MV2, which is the **more subtle** of the two traces.



Selected input: record  
ltstdb/s20101 , annotator atr ,  
from [10:25:40.000 01/11/1994]  
to [10:25:50.000 01/11/1994]

Long Term ST Database (ltstdb)

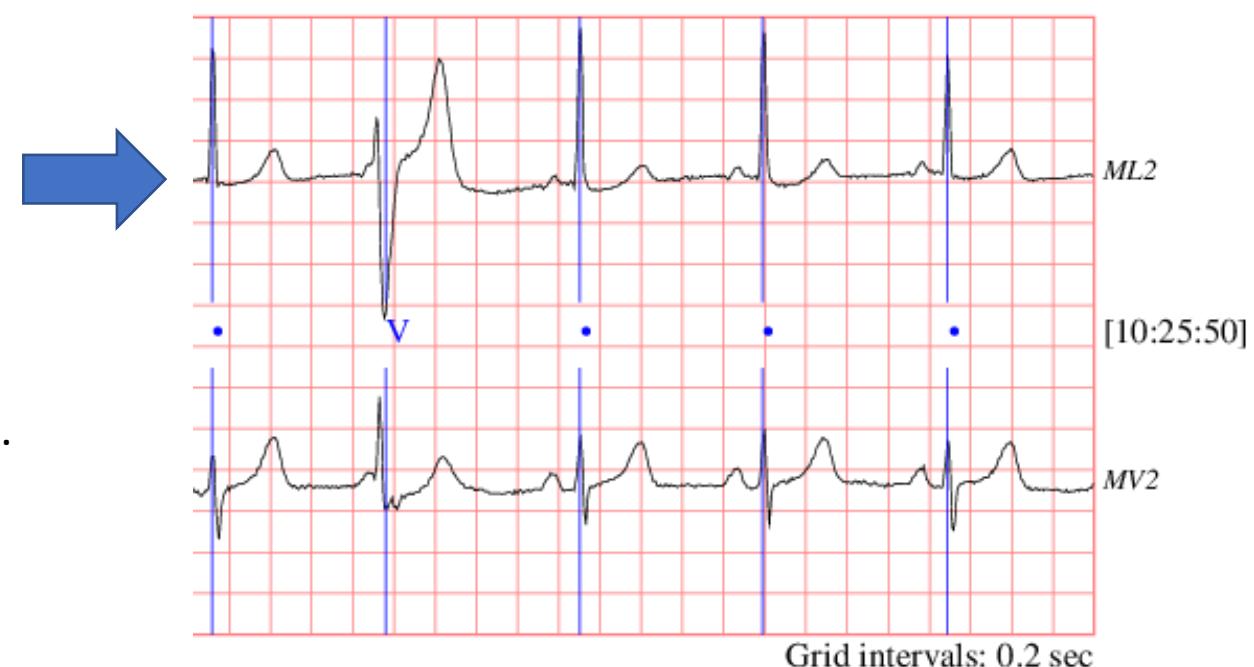
```
pschart -a atr -c "" -C -E -G -CG 1 .5 .5 -Cs 0  
0 0 -H -I -P 300x250 -m 20 20 5 5 -M -n 0 \ -S  
4 2 -t 25 -T "" -v 10 -w 0.5 script >chart.ps  
convert -density 100x100 chart.ps chart.png
```



# UCR\_Anomaly\_s20101mML2\_12000\_35774\_35874.txt

This dataset contains normal beats, with a single ventricular beat between 35774 and 35874.

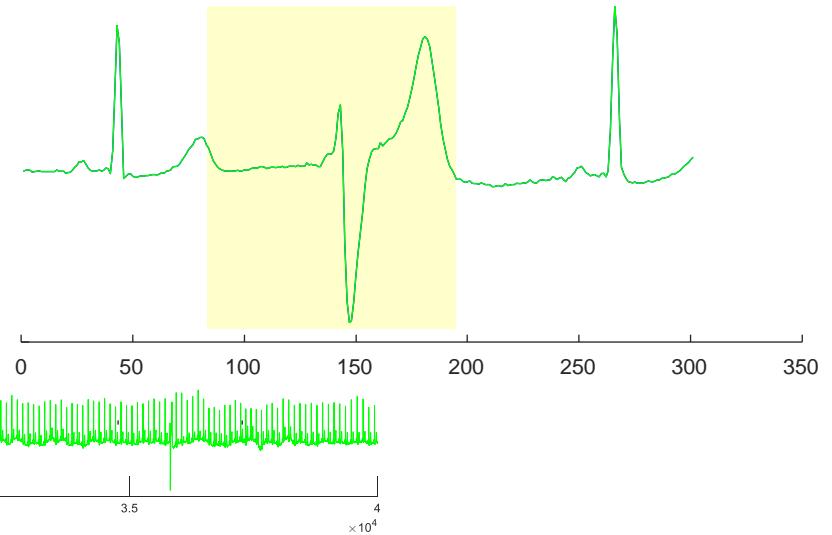
Here we use ML2, which is the **more obvious** of the two traces.



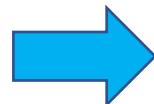
Selected input: record  
ltstdb/s20101 , annotator atr ,  
from [10:25:40.000 01/11/1994]  
to [10:25:50.000 01/11/1994]

Long Term ST Database (ltstdb)

[pschart](#) -a atr -c "" -C -E -G -CG 1 .5 .5 -Cs 0  
0 0 -H -I -P 300x250 -m 20 20 5 5 -M -n 0 \ -S  
4 2 -t 25 -T "" -v 10 -w 0.5 [script >chart.ps](#)  
[convert](#) -density 100x100 chart.ps [chart.png](#)



# UCR\_Anomaly\_ltstdbs30791ES\_20000\_52600\_52800.txt

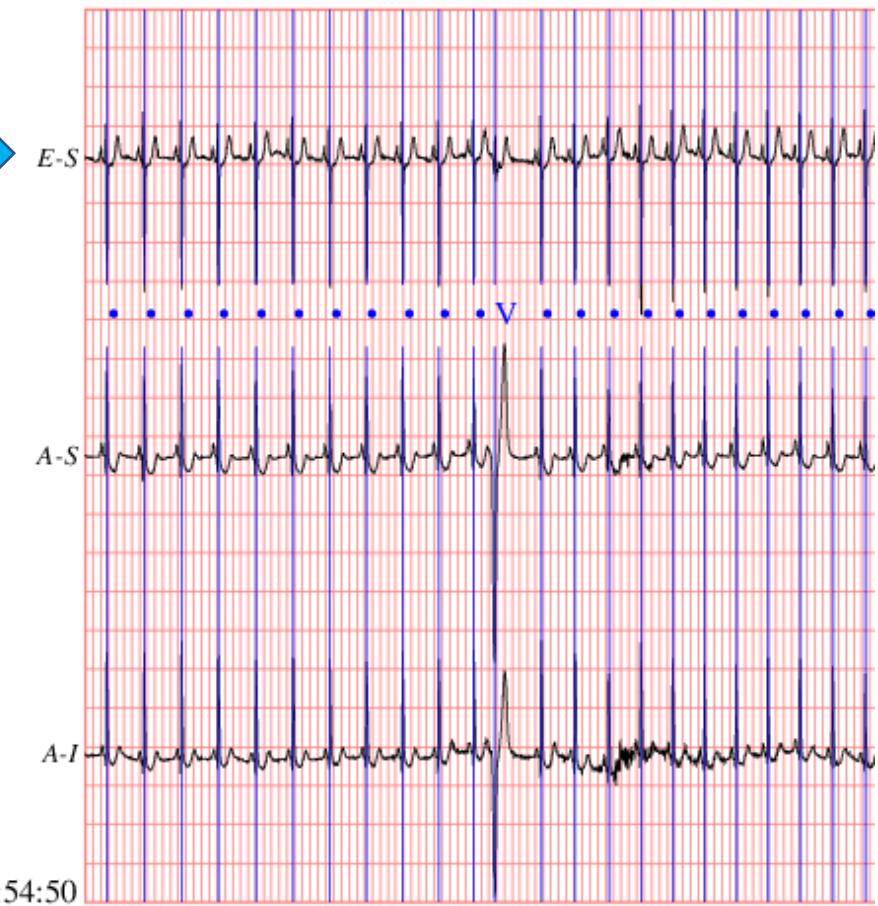
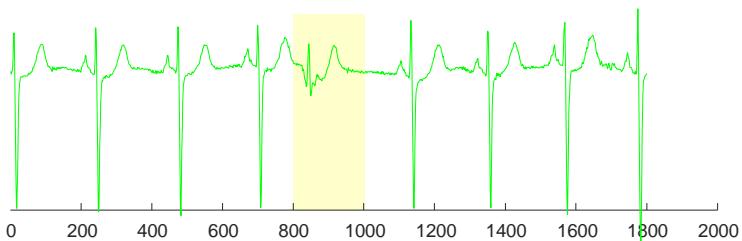


This dataset contains normal beats, with a single ventricular beat between 52600 and 52800.

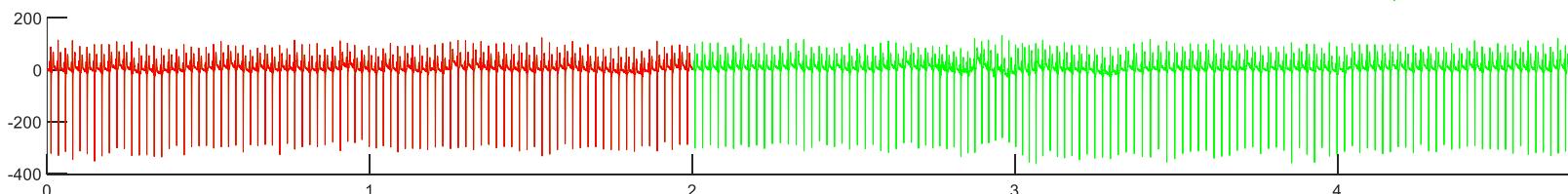
Here we use E-S, which is the **more subtle** of the three traces.

There are some noisy regions in the test set, which some algorithms might consider anomalies. However, there are just as noisy regions in the training set.

Likewise, there is slight wandering baseline in the test data, but we also see this in the train data.



```
pschart -a atr -c "" -C -E -G -CG 1 .5 .5 -Cs 0 0 0 -H -I -P  
360x250 -m 20 20 5 5 -M -n 0 \ -S 4 2 -t 5 -T "" -v 10 -w 0.5  
script >chart.ps  
convert -density 100x100 chart.ps chart.png
```



Selected  
input: record  
ltstdb/s30791 ,  
annotator atr ,  
from 8:54:50.000  
to 8:55:50.000

Long Term ST  
Database (ltstdb)

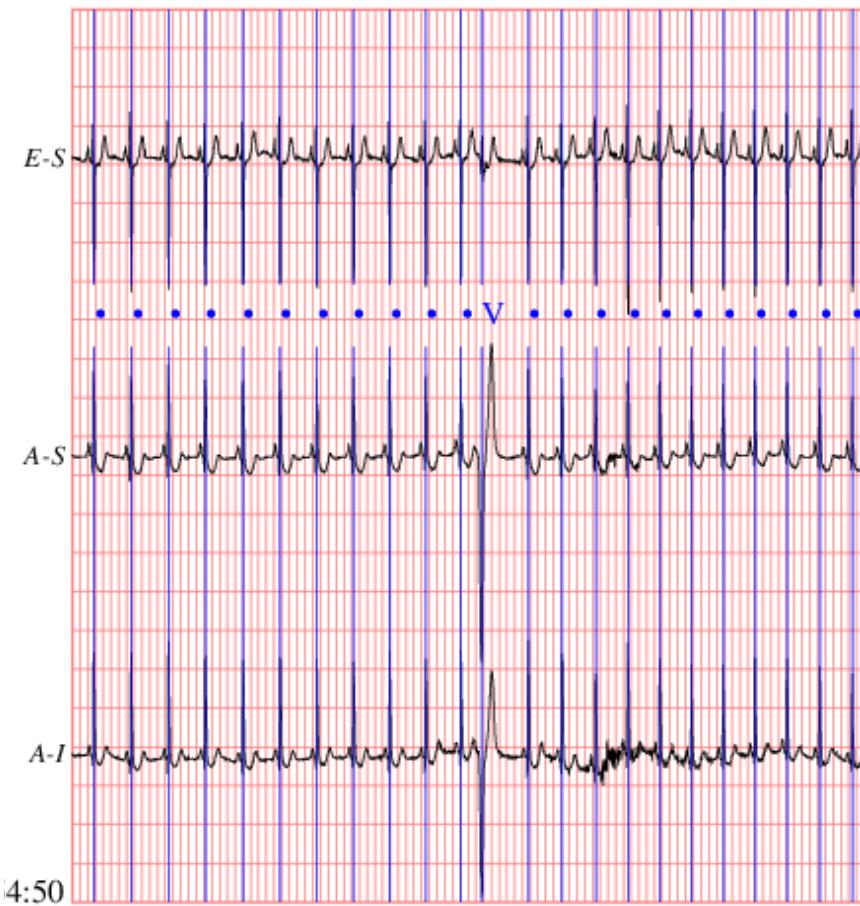
# UCR\_Anomaly\_ltstdbs30791AS\_23000\_52600\_52800.txt

This dataset contains normal beats, with a single ventricular beat between 52600 and 52800.

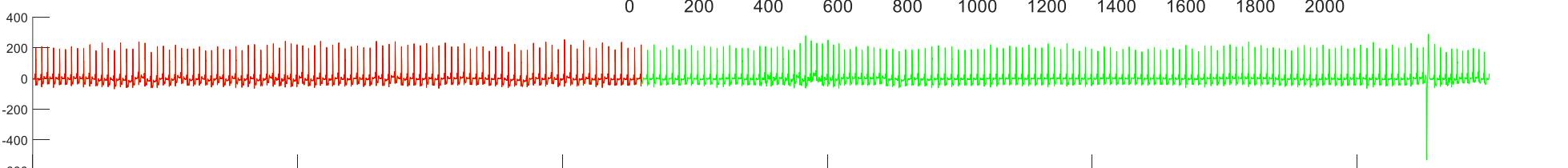
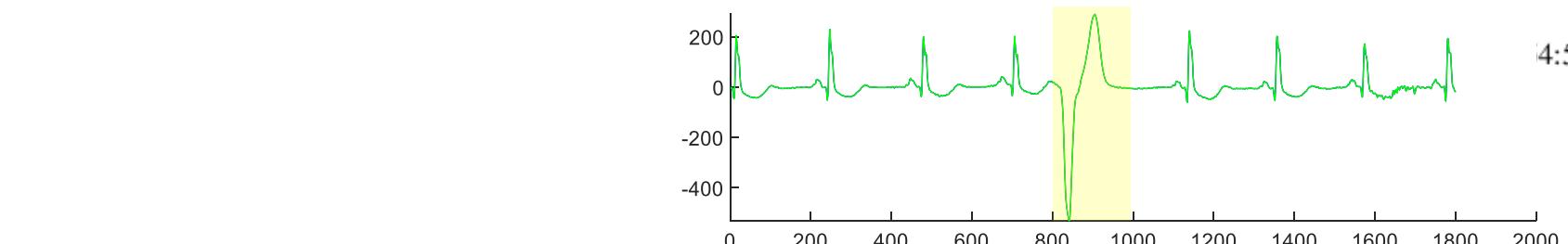
Here we use A-S, which is the **least subtle** of the three traces.

There are some noisy regions in the test set, which some algorithms might consider anomalies. However, there are just as noisy regions in the training set.

Likewise, there is slight wandering baseline in the test data, but we also see this in the train data.



```
pschart -a atr -c "" -C -E -G -CG 1 .5 .5 -Cs 0 0 0 -H -I -P  
360x250 -m 20 20 5 5 -M -n 0 \ -S 4 2 -t 5 -T "" -v 10 -w 0.5  
script >chart.ps  
convert -density 100x100 chart.ps chart.png
```



Selected  
input: record  
ltstdb/s30791 ,  
annotator atr ,  
from 8:54:50.000  
to 8:55:50.000

Long Term ST  
Database (ltstdb)

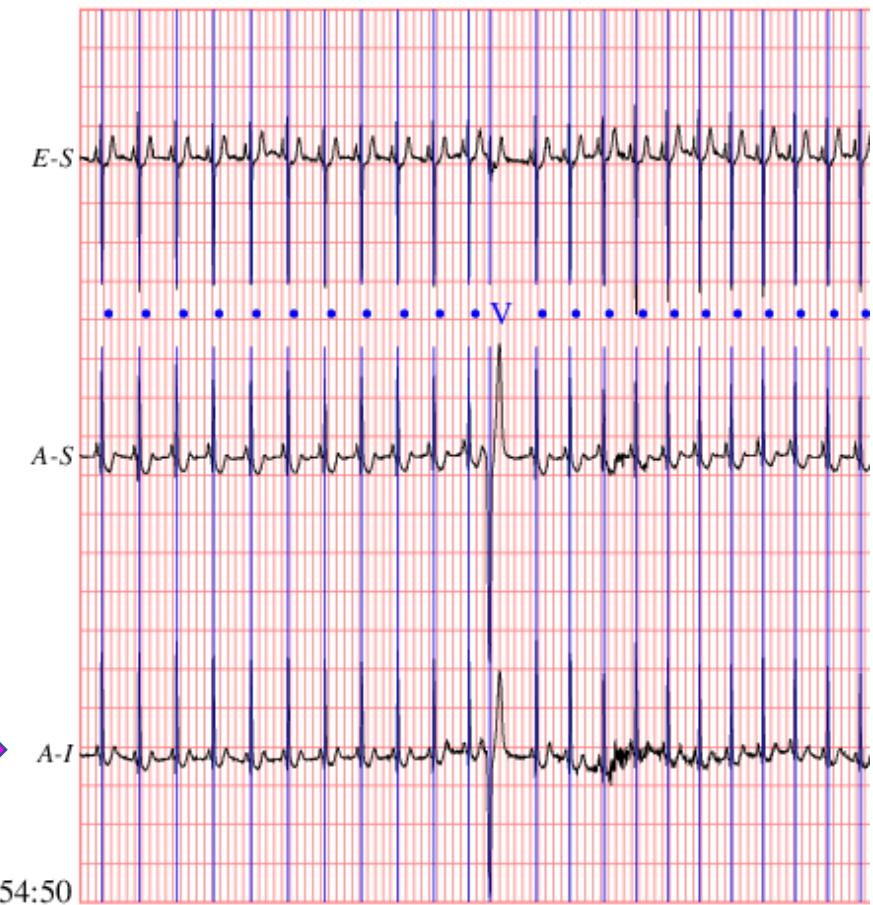
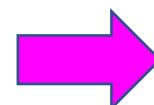
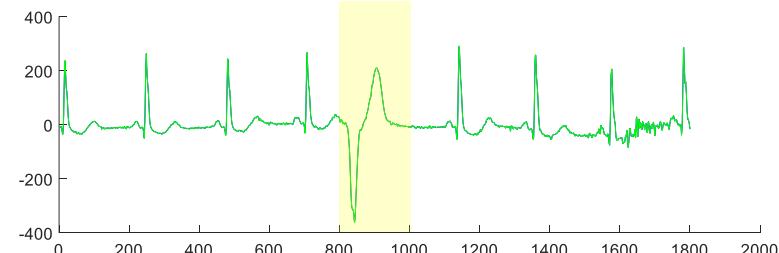
# UCR\_Anomaly\_ltstdbs30791AI\_17555\_52600\_52800.txt

This dataset contains normal beats, with a single ventricular beat between 52600 and 52800.

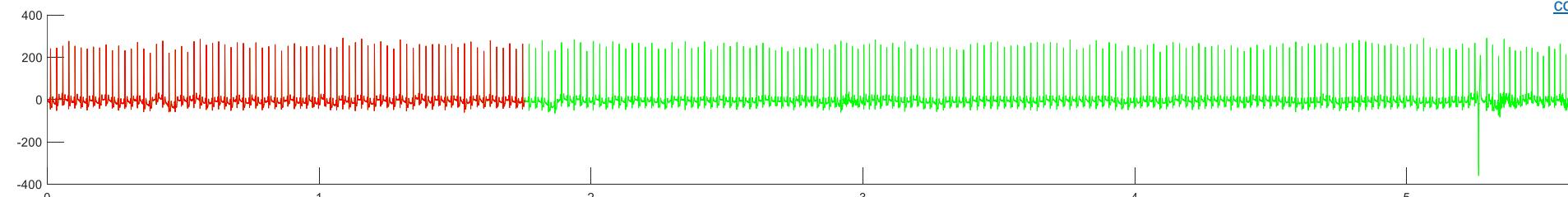
Here we use A-I.

There are some noisy regions in the test set, which some algorithms might consider anomalies. However, there are just as noisy regions in the training set.

Likewise, there is slight wandering baseline in the test data, but we also see this in the train data.



```
pschart -a atr -c "" -C -E -G -CG 1 .5 .5 -Cs 0 0 0 -H -I -P  
360x250 -m 20 20 5 5 -M -n 0 \ -S 4 2 -t 5 -T "" -v 10 -w 0.5  
script >chart.ps  
convert -density 100x100 chart.ps chart.png
```



Selected  
input: record  
ltstdb/s30791 ,  
annotator atr ,  
from 8:54:50.000  
to 8:55:50.000

Long Term ST  
Database (ltstdb)

# UCR\_Anomaly\_qtdbSel100MLII\_4000\_13400\_13800.txt

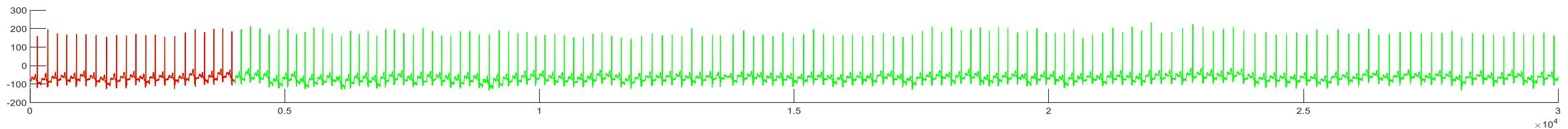
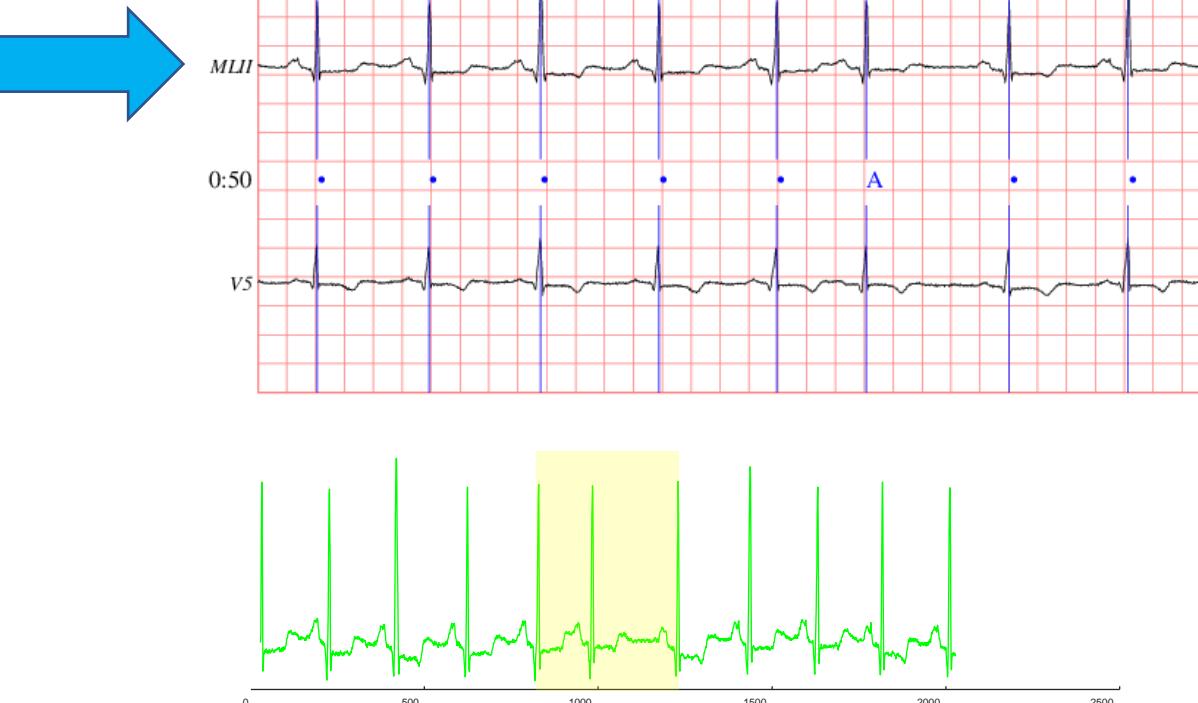
```
pschart -a atr -c "" -C -E -G -CG 1 .5 .5 -Cs 0 0 0 -H -I -P  
300x250 -m 20 20 5 5 -M -n 0 \ -S 4 2 -t 25 -T "" -v 10 -w  
0.5 script >chart.ps  
convert -density 100x100 chart.ps chart.png
```

**Selected**  
input: record  
qtdb/sel100 ,  
annotator atr , from  
0:50.000 to  
1:00.000

[QT Database \(qtdb\)](#)

This dataset has a single bad heartbeat, an Atrial premature beat between 13400 and 13800.  
By most standards, this is a subtle anomaly.

This is the [MLII](#) trace



# UCR\_Anomaly\_qtdbSel1005V\_4000\_12400\_12800.txt

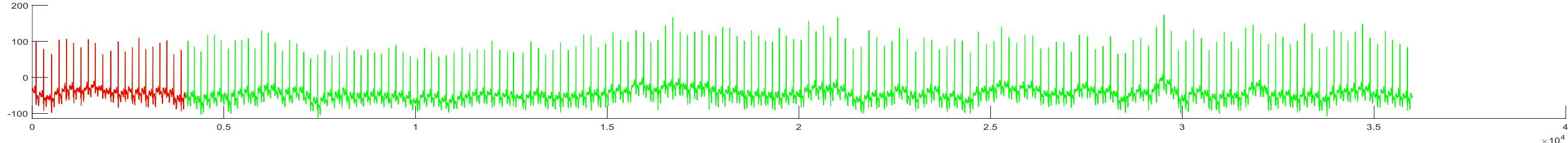
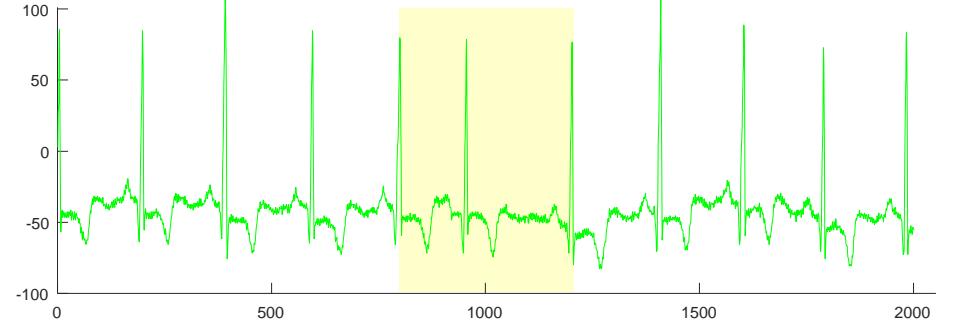
```
pschart -a atr -c "" -C -E -G -CG 1 .5 .5 -Cs 0 0 0 -H -I -P  
300x250 -m 20 20 5 5 -M -n 0 \ -S 4 2 -t 25 -T "" -v 10 -w  
0.5 script >chart.ps  
convert -density 100x100 chart.ps chart.png
```

**Selected**  
**input:** record  
qtdb/sel100 ,  
annotator atr , from  
0:50.000 to  
1:00.000

[QT Database \(qtdb\)](#)

This dataset has a single bad heartbeat, an Atrial premature beat between 12400 and 12800.  
By most standards, this is a subtle anomaly.

This is the **V5** trace. It is shifted by 1000 datapoints from  
UCR\_Anomaly\_qtdbSel100MLII\_4000\_13400\_13800.txt



```
wfdb2mat -r qtdb/sel840 -f  
10 -t 899.996 -l s1000000  
>sel840m.info
```

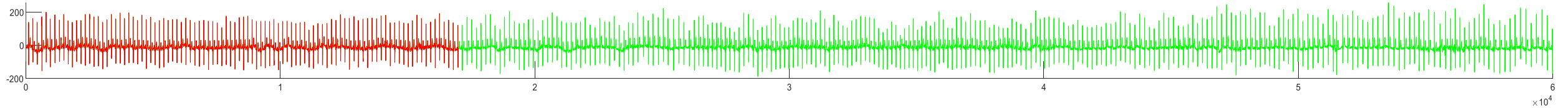
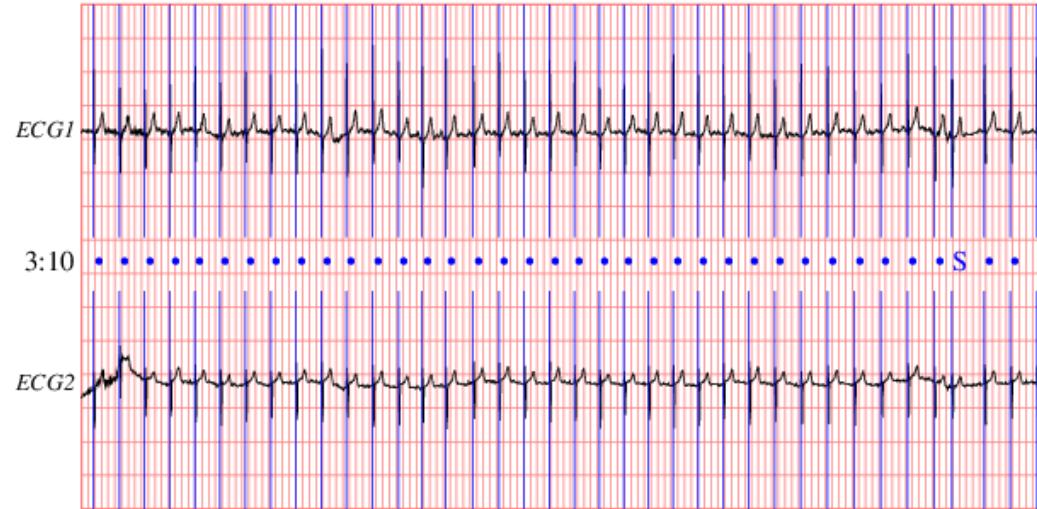
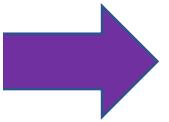
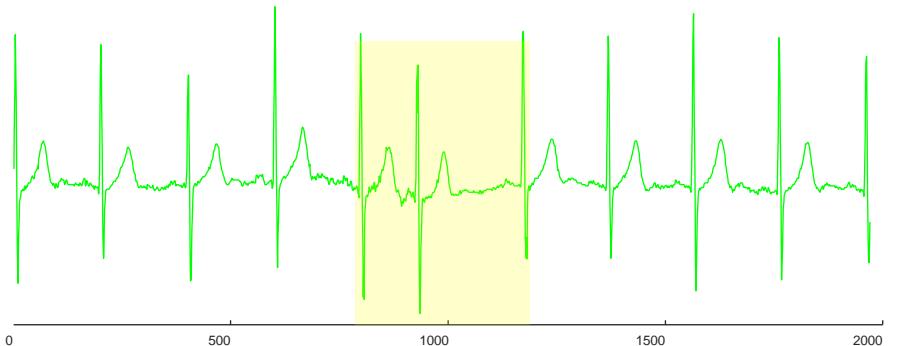
# UCR\_Anomaly\_sel840mECG1\_17000\_51370\_51740.txt

Selected  
input: record  
qtdb/sel840 ,  
annotator atr ,  
from 0:10.000 to  
14:59.996

This dataset has a single bad heartbeat, a supraventricular beat between 51370 and 51740.

This is **ECG1**

By most standards, this is a subtle anomaly.



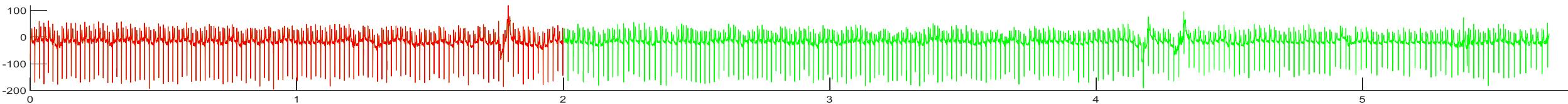
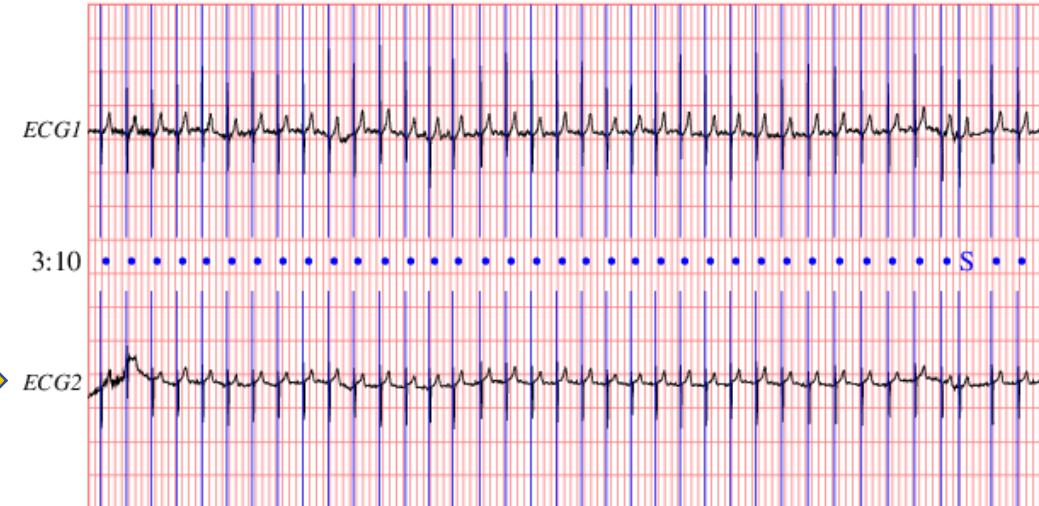
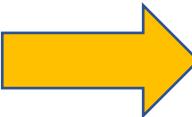
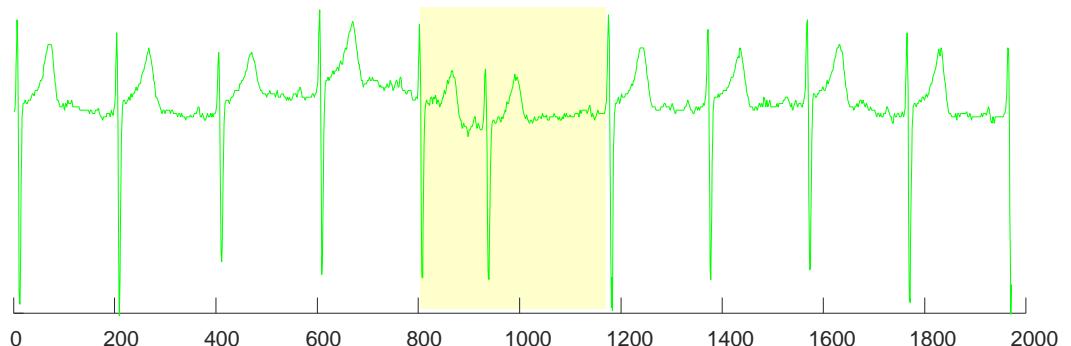
```
wfdb2mat -r qtdb/sel840 -f  
10 -t 899.996 -l s1000000  
>sel840m.info
```

# UCR\_Anomaly\_sel840mECG2\_20000\_49370\_49740.txt

This dataset has a single bad heartbeat, a supraventricular beat between 49370 and 49740.

This is ECG2

This is shifted by 2000 datapoints from [UCR\\_Anomaly\\_sel840mECG1\\_17000\\_51370\\_51740.txt](#)  
By most standards, this is a subtle anomaly.

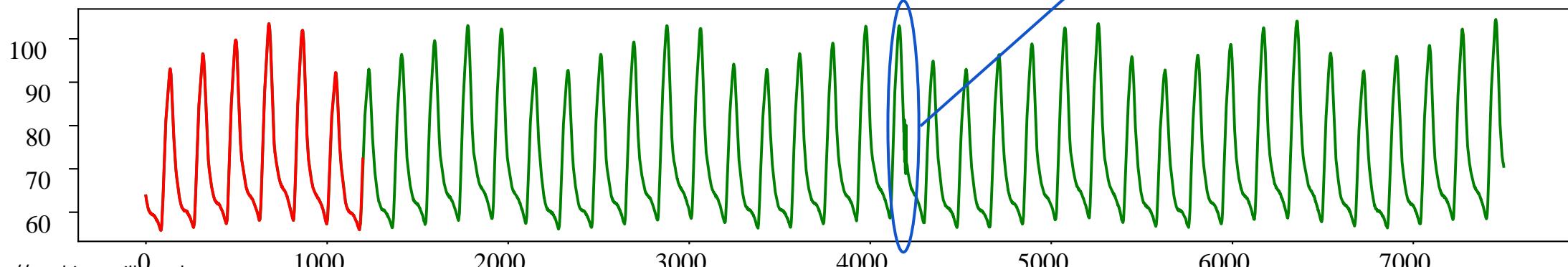
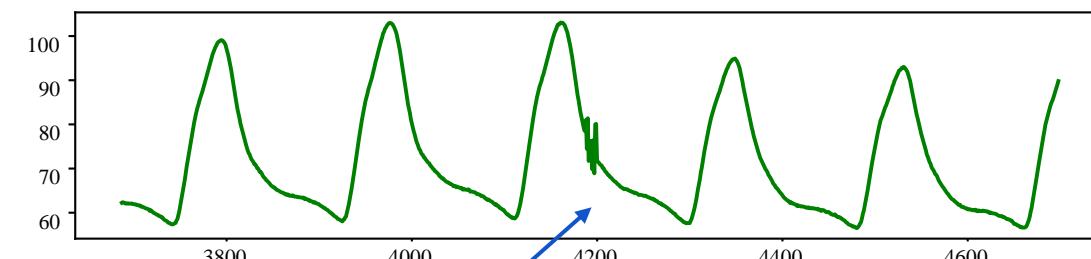
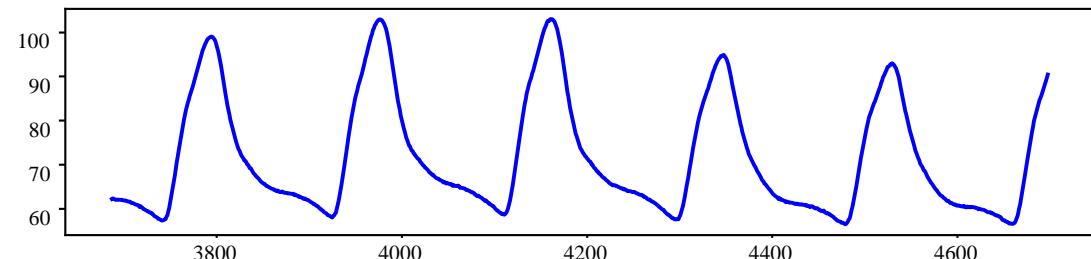


# UCR\_Anomaly\_InternalBleeding16\_1200\_4187\_4199.txt

The data comes from an internal bleeding dataset<sup>[1]</sup>. From the dataset, we extracted the arterial blood pressure measurements of pigs.

The anomaly is synthetic. We generated random numbers to add noise from 4187 to 4199.

Blue is original data, green is data after anomaly was introduced

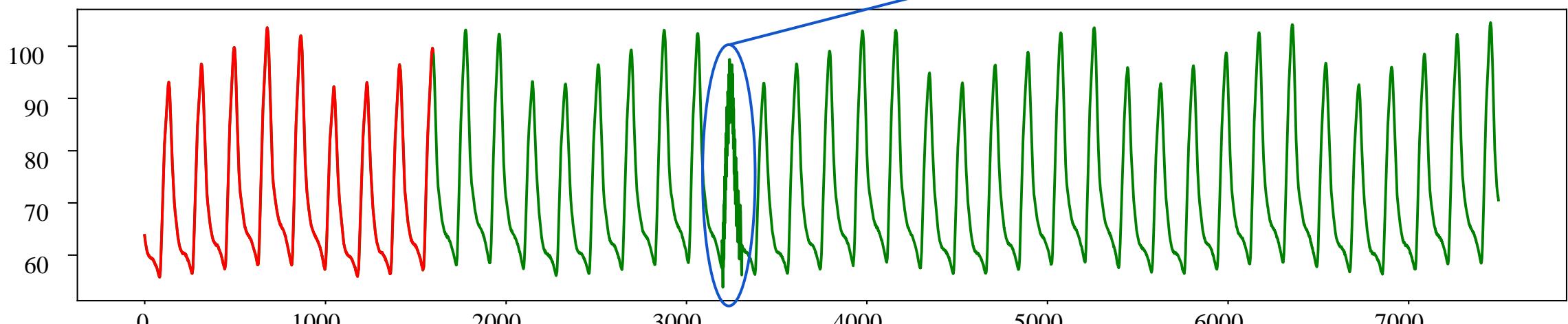
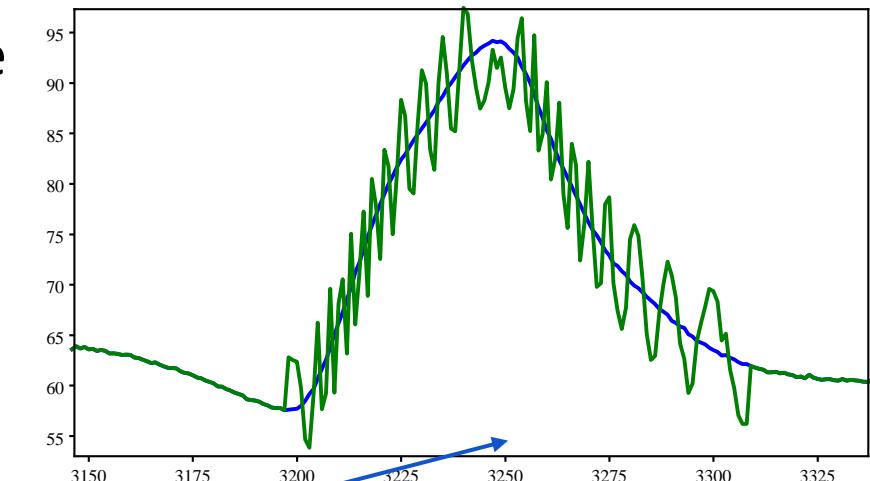


# UCR\_Anomaly\_InternalBleeding17\_1600\_3198\_3309.txt

The data comes from an internal bleeding dataset<sup>[1]</sup>. From the dataset, we extracted the arterial blood pressure measurements of pigs.

The anomaly is synthetic. We added a series of sine waves in one cycle to make a piece of smooth data become fluctuate and rough.

Blue is original data, green is data after anomaly was introduced



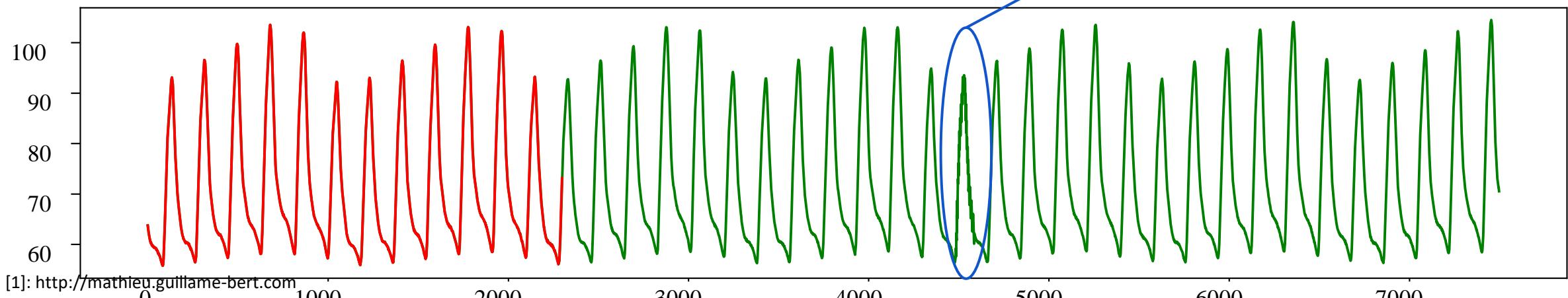
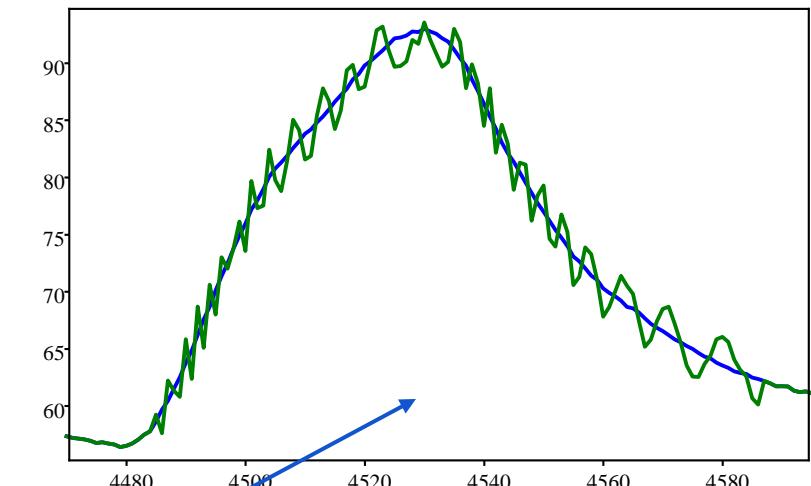
[1]: <http://mathieu.guillame-bert.com>

# UCR\_Anomaly\_InternalBleeding18\_2300\_4485\_4587.txt

The data comes from an internal bleeding dataset<sup>[1]</sup>. From the dataset, we extracted the arterial blood pressure measurements of pigs.

Blue is original data, green is data after anomaly was introduced

The anomaly is synthetic. We added a series of sine waves in one cycle to make a piece of smooth data become fluctuate and rough. This is a more subtle version of UCR\_Anomaly\_InternalBleeding16\_1200\_4187\_4199.txt.



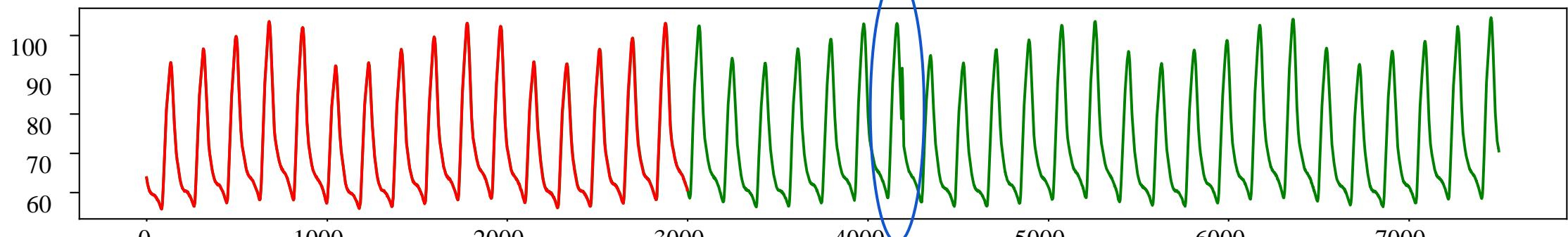
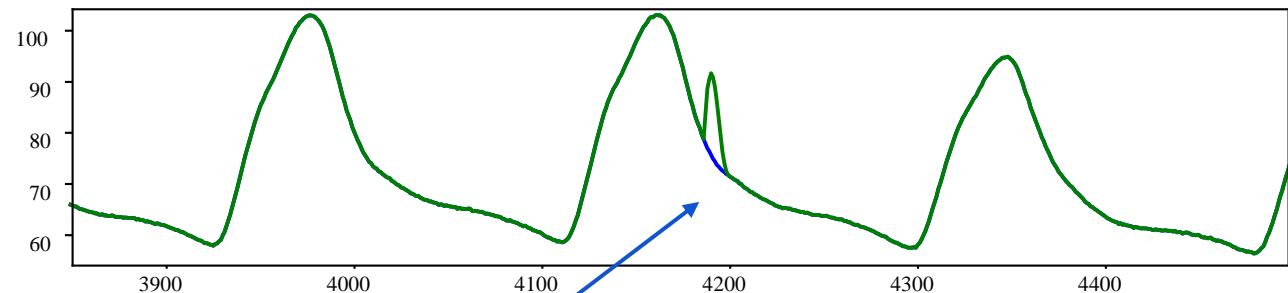
[1]: <http://mathieu.guillame-bert.com>

# UCR\_Anomaly\_InternalBleeding19\_3000\_4187\_4197.txt

The data comes from an internal bleeding dataset<sup>[1]</sup>. From the dataset, we extracted the arterial blood pressure measurements of pigs.

The anomaly is synthetic. We added a secondary peak to a cycle of data that has only one peak.

Blue is original data, green is data after anomaly was introduced



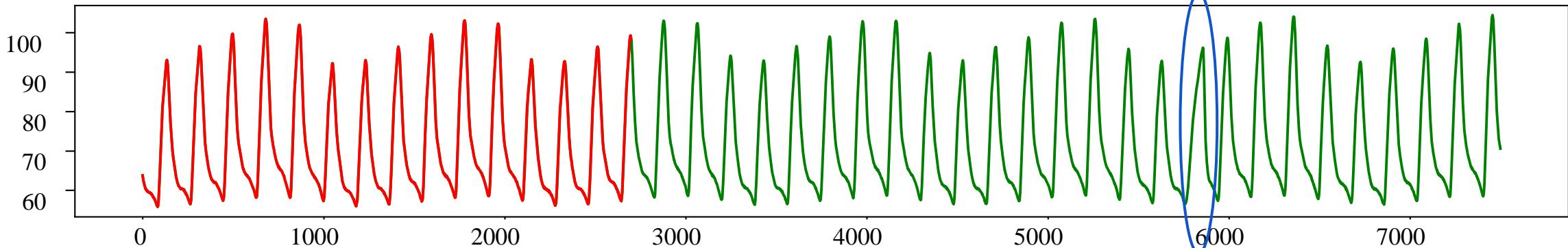
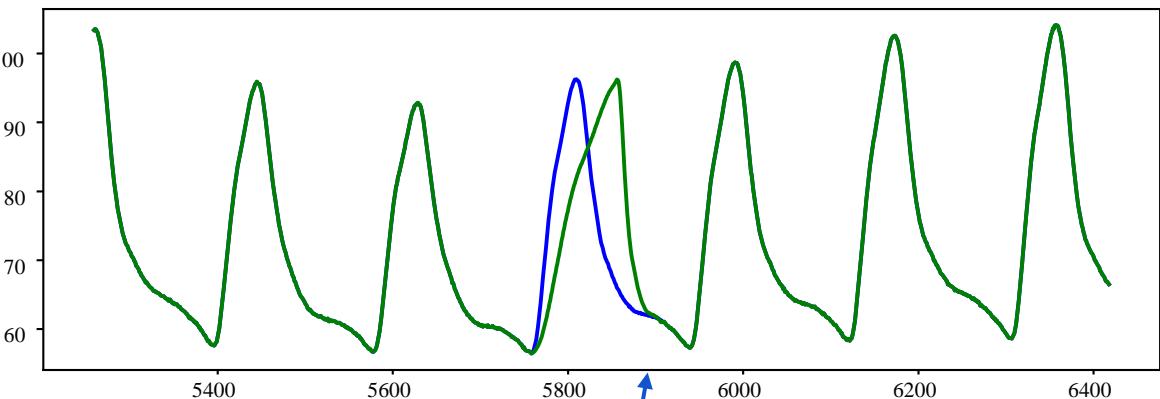
[1]: <http://mathieu.guillame-bert.com>

# UCR\_Anomaly\_InternalBleeding20\_2700\_5759\_5919.txt

The data comes from an internal bleeding dataset<sup>[1]</sup>. From the dataset, we extracted the arterial blood pressure measurements of pigs.

The anomaly is synthetic. We performed resampling (upsampling + downsampling) operations to skew the data in one cycle.

Blue is original data, green is data after anomaly was introduced

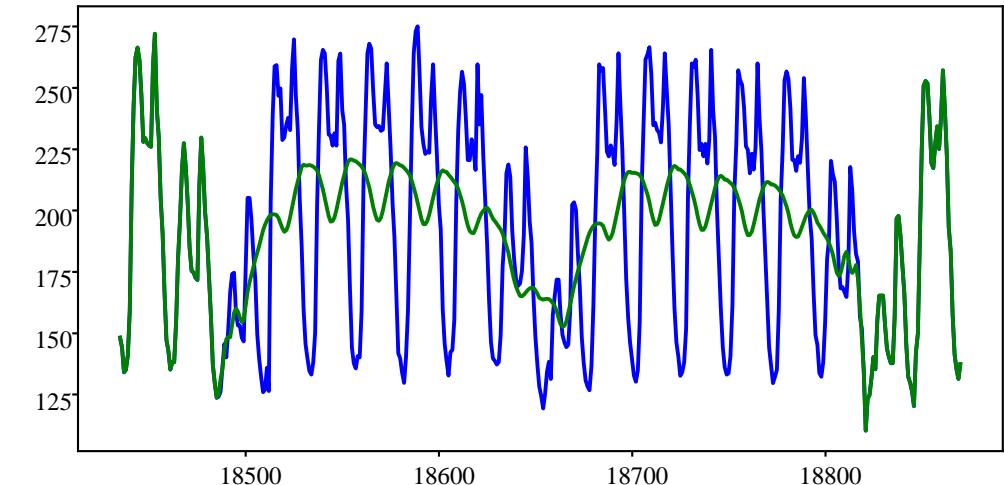


[1]: <http://mathieu.guillame-bert.com>

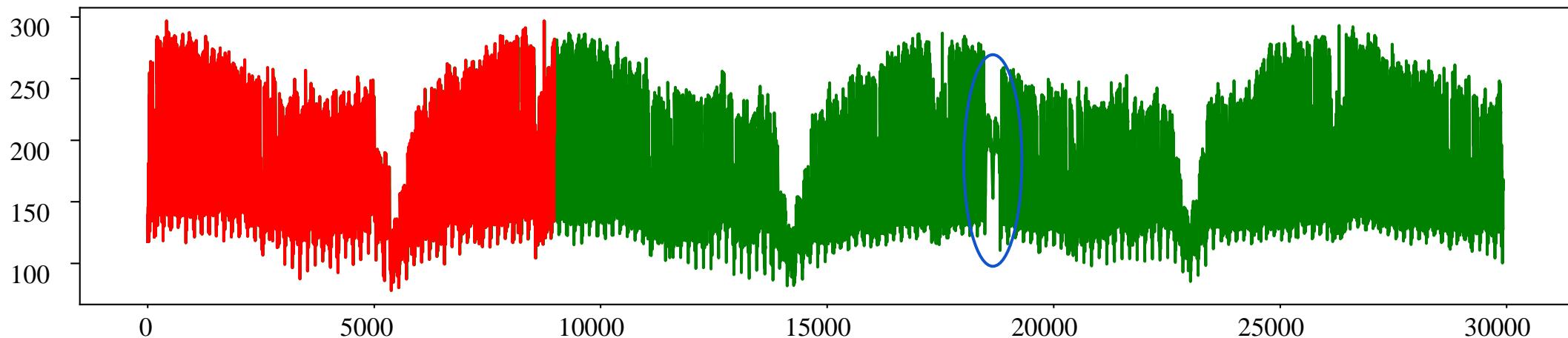
# UCR\_Anomaly\_PowerDemand1\_9000\_18485\_18821.txt

The data comes from an Italian power demand dataset. (from 1/1/1995 to 5/31/1998)

The anomaly is synthetic. We ran the moving average algorithm on random two-week power supply data (2/9/1997 3:00 - 2/23/1997 3:00).



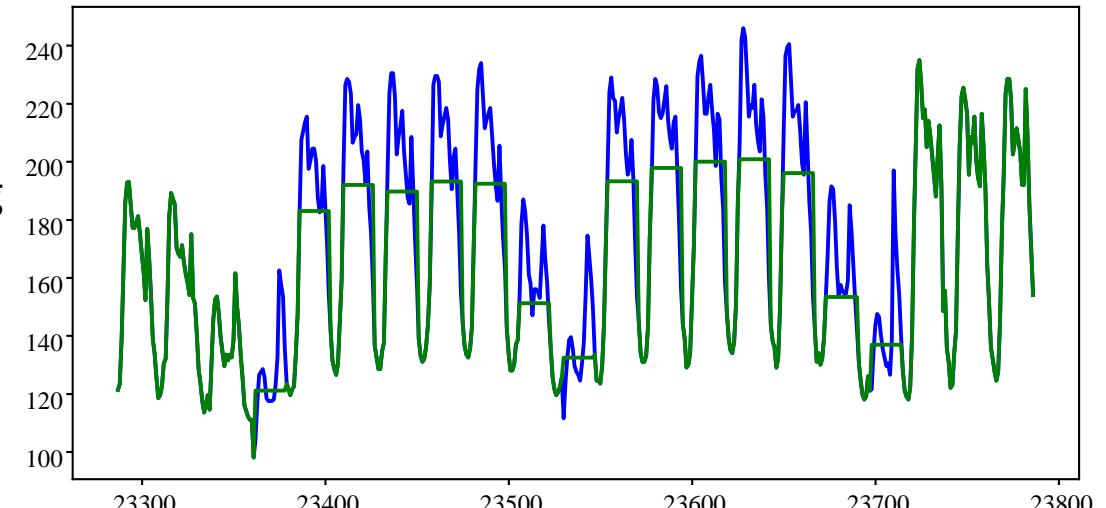
Blue is original data, green is data after anomaly was introduced



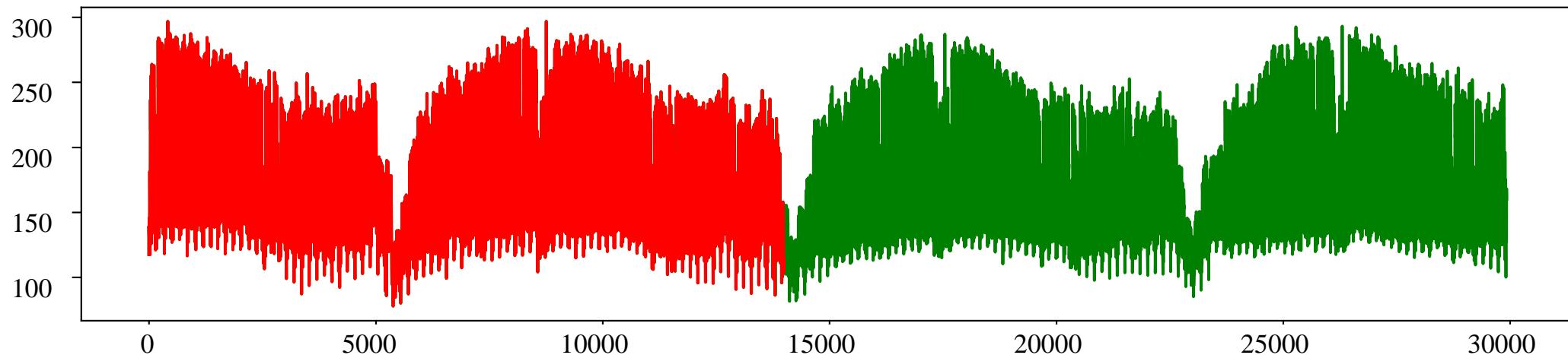
# UCR\_Anomaly\_PowerDemand2\_14000\_23357\_23717.txt

The data comes from an Italian power demand dataset. (from 1/1/1995 to 5/31/1998)

The anomaly is synthetic. We created an anomaly by removing peaks and adding constant regions for half a month's power supply data. (8/31/1997 - 9/15/1997)



Blue is original data, green is data after anomaly was introduced

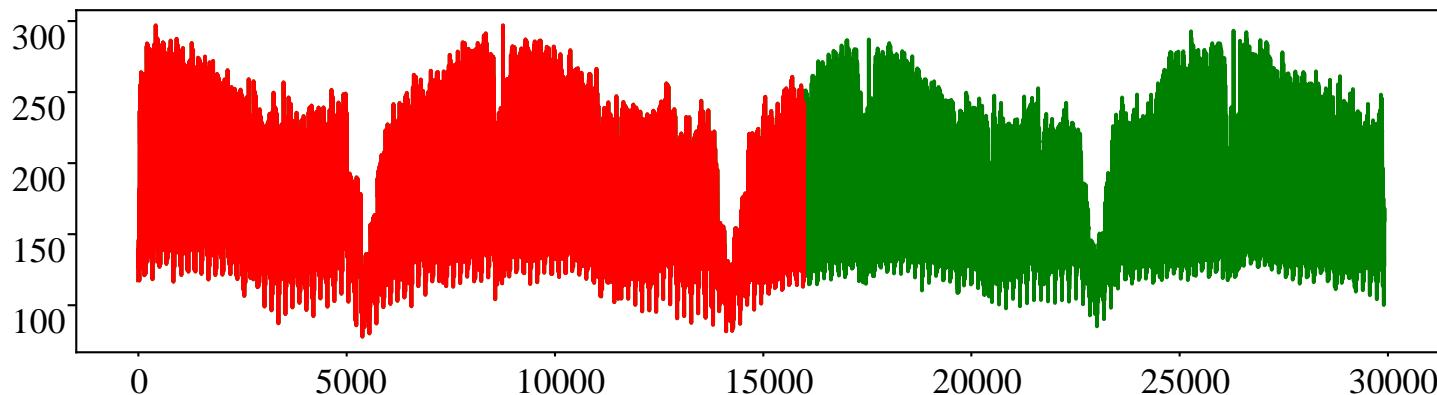
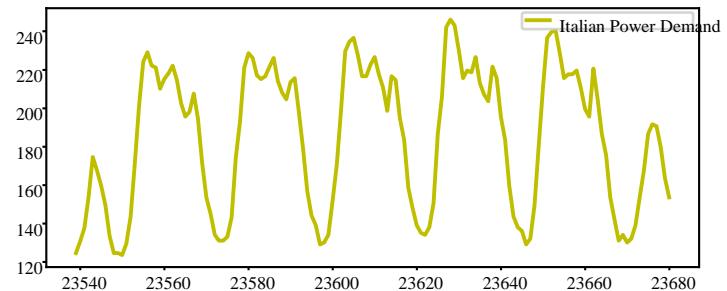
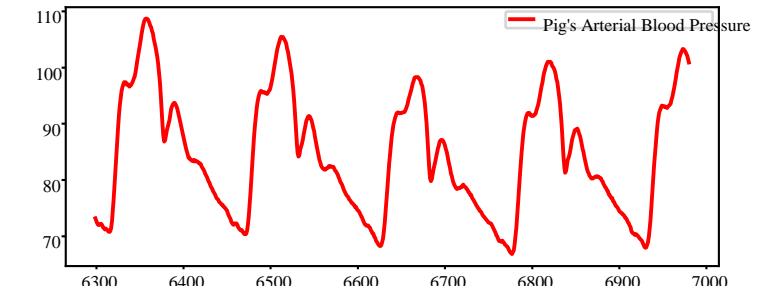


# UCR\_Anomaly\_PowerDemand3\_16000\_23405\_23477.txt

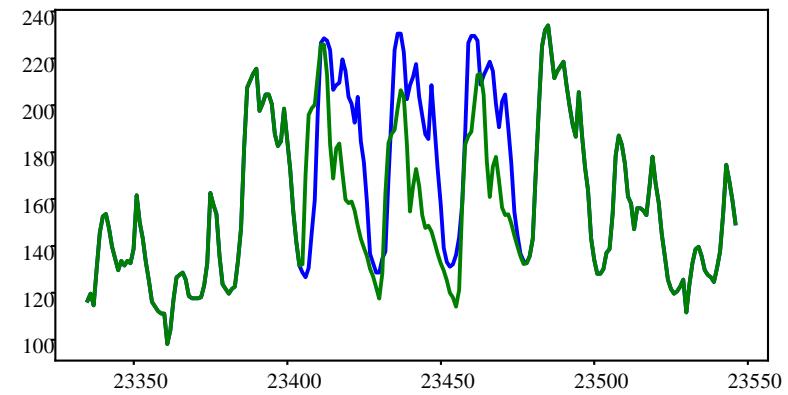
The data comes from an Italian power demand dataset.

(from 1/1/1995 to 5/31/1998) **The red curve indicates the pig arterial blood pressure, and the yellow curve indicates the Italian power demand. These two curves look similar.**

The anomaly is synthetic. We replaced random three-day power supply data with a segment of arterial blood pressure data. ( The X of arterial blood pressure data is scaled by resampling and the Y is adjusted by a level shift.)

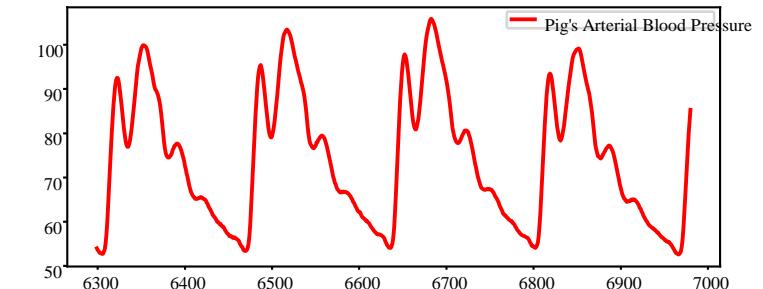


Blue is original data, green is data after anomaly was introduced

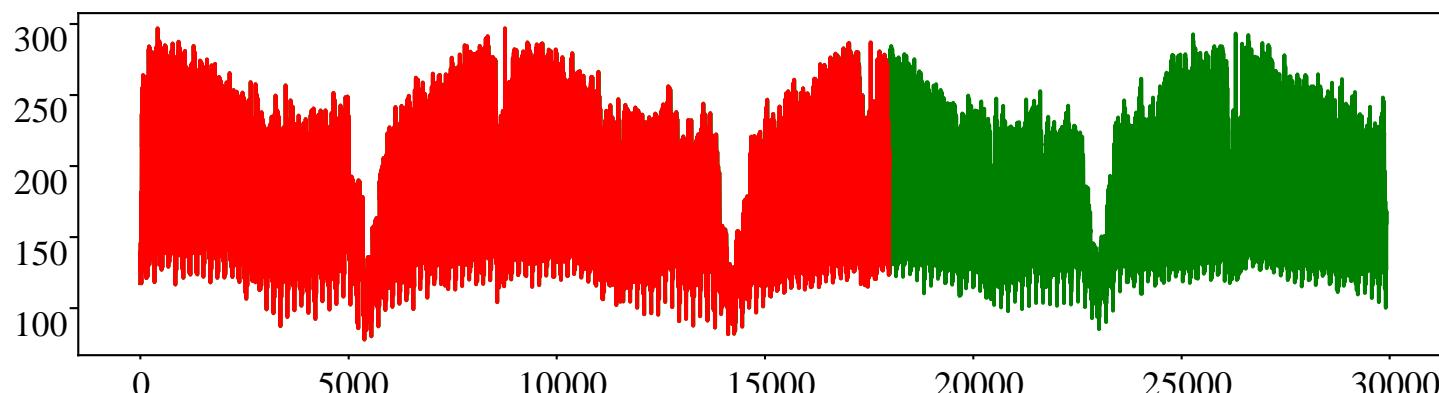
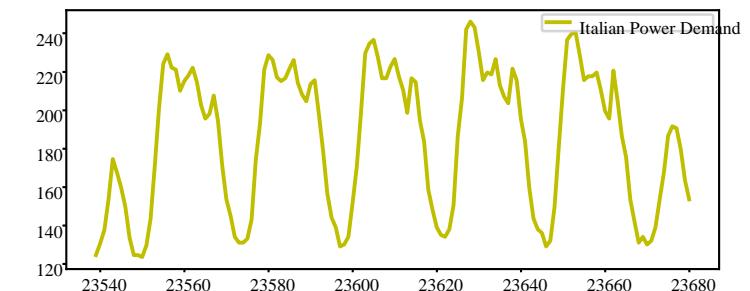


# UCR\_Anomaly\_PowerDemand4\_18000\_24005\_24077.txt

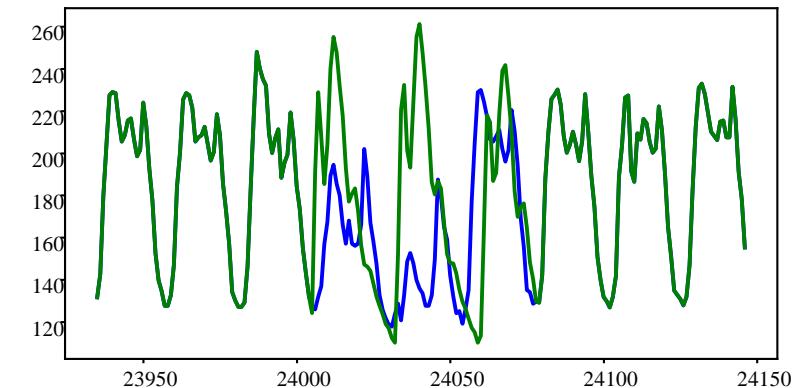
The data comes from an Italian power demand dataset. (from 1/1/1995 to 5/31/1998) **The red curve indicates the pig arterial blood pressure, and the yellow curve indicates the Italian power demand.** These two curves look similar.



The anomaly is synthetic. We replaced random three-day power supply data with a segment of arterial blood pressure data. ( The X of arterial blood pressure data is scaled by resampling and the Y is adjusted by a level shift.)



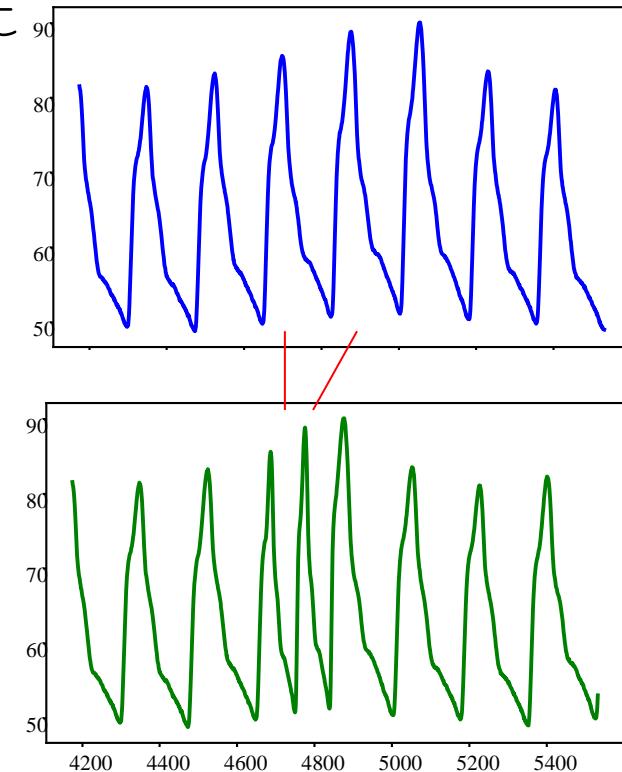
Blue is original data, green is data after anomaly was introduced



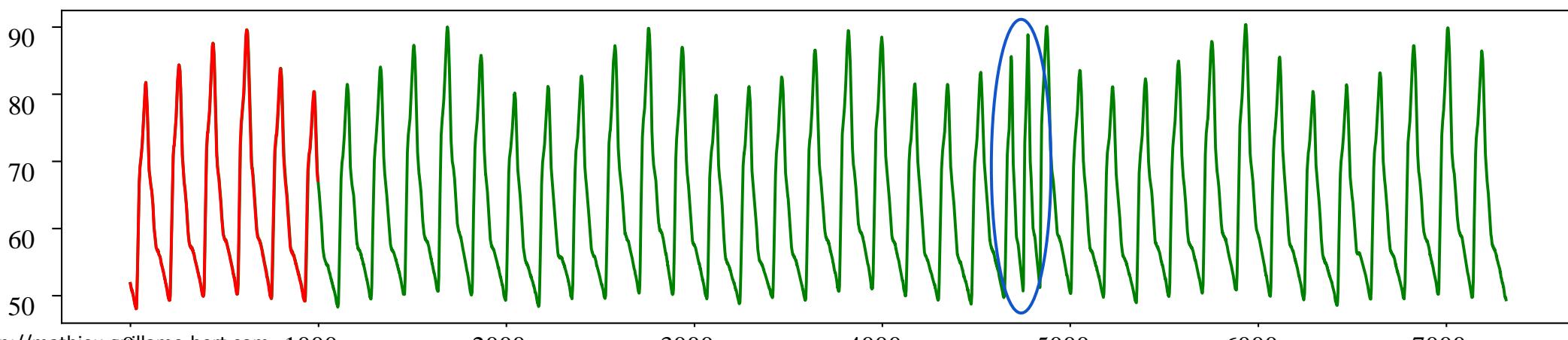
# UCR\_Anomaly\_InternalBleeding4\_1000\_4675\_5033.txt

The data comes from an internal bleeding dataset<sup>[1]</sup>. From the dataset, we extracted the arterial blood pressure measurements of pigs.

The anomaly is synthetic. We performed a downsampling operation on a short fraction of normal arterial blood pressure data. From 4675 to 5033, the values at every two time points are combined by averaging.



Blue is original data, green is data after  
anomaly was introduced

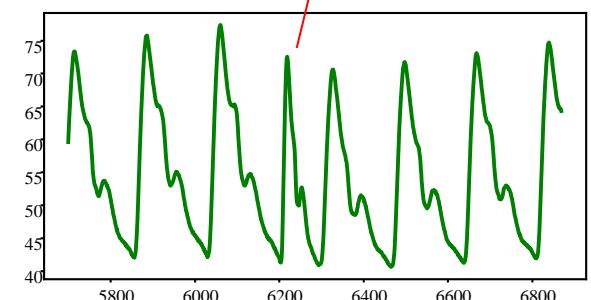
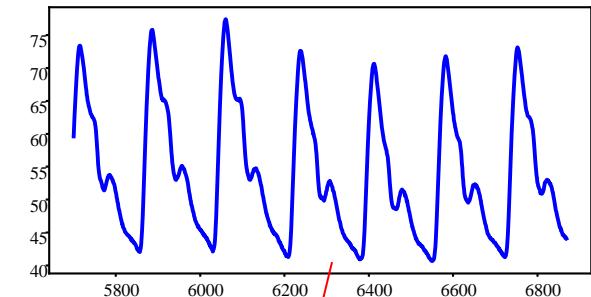


## UCR\_Anomaly\_InternalBleeding5\_4000\_6200\_6370.txt

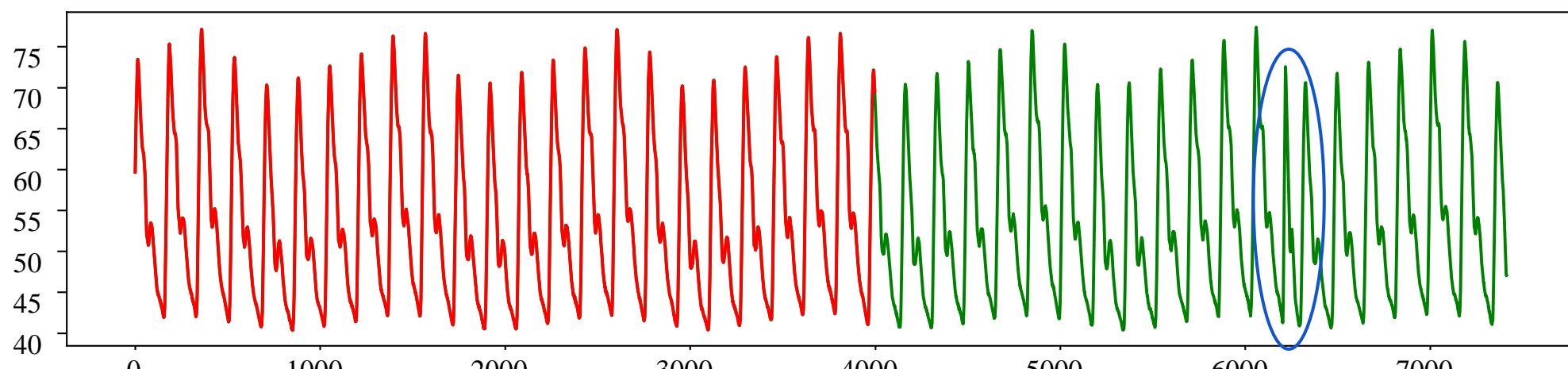
The data comes from an internal bleeding dataset<sup>[1]</sup>. From the dataset, we extracted the arterial blood pressure measurements of pigs.

The anomaly is synthetic. We performed a downsampling operation on a short fraction of normal arterial blood pressure data. From 6200 to 6370, the values at every two time points are combined by averaging.

This is a more subtle version of UCR\_Anomaly\_InternalBleeding4\_1000\_4675\_5033.txt



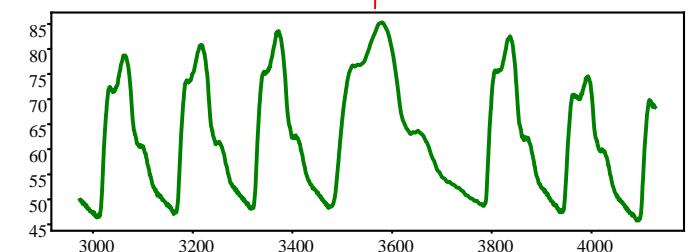
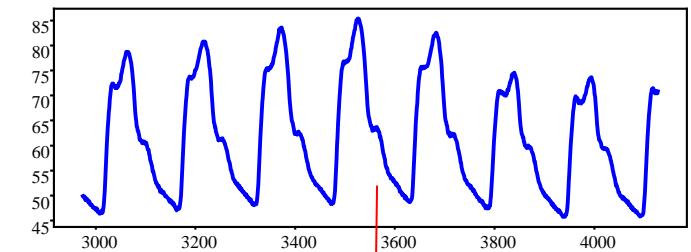
Blue is original data, green is data after anomaly was introduced



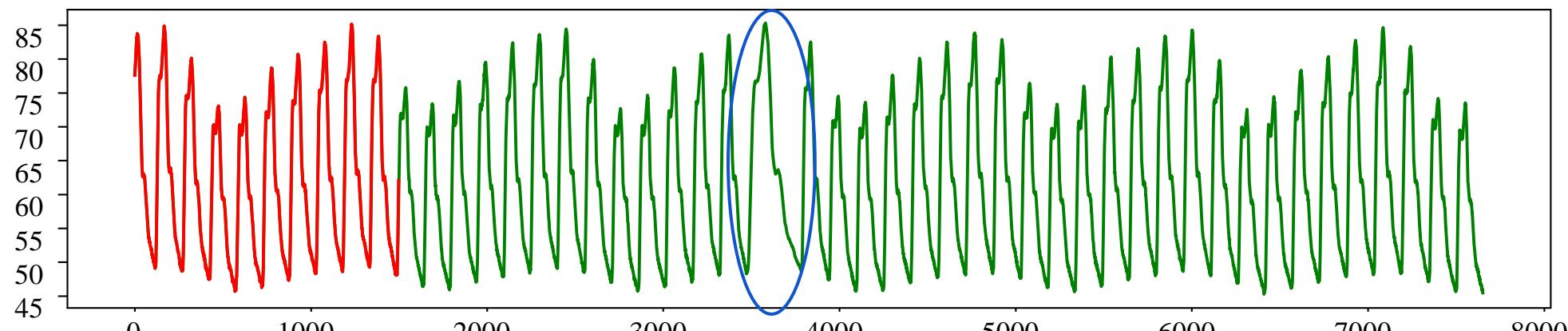
# UCR\_Anomaly\_InternalBleeding6\_1500\_3474\_3629.txt

The data comes from an internal bleeding dataset<sup>[1]</sup>. From the dataset, we extracted the arterial blood pressure measurements of pigs.

The anomaly is synthetic. We performed an upsampling operation on a short fraction of normal arterial blood pressure data. From 3393 to 3570, the frequency of the time series is doubled, and the missing data is filled with the value of the previous time point.



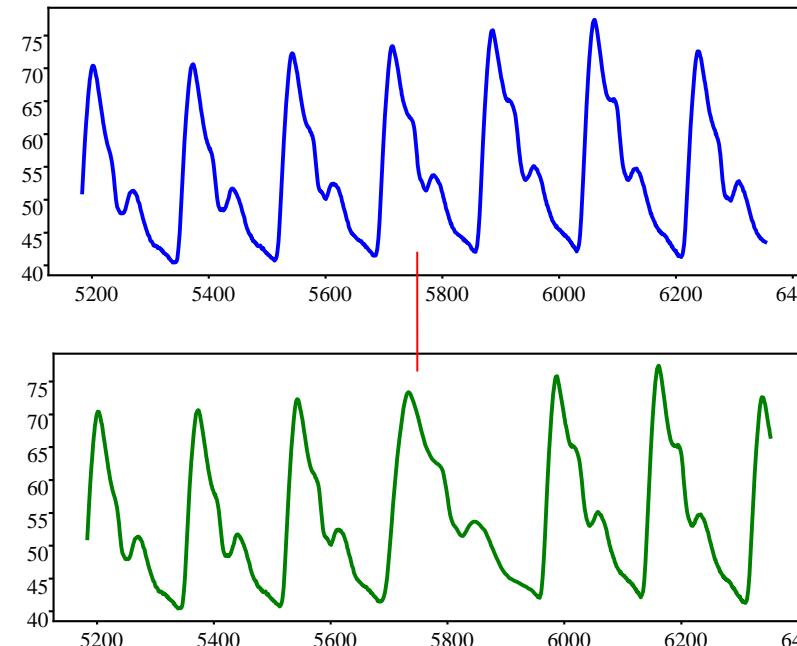
Blue is original data, green is data after anomaly was introduced



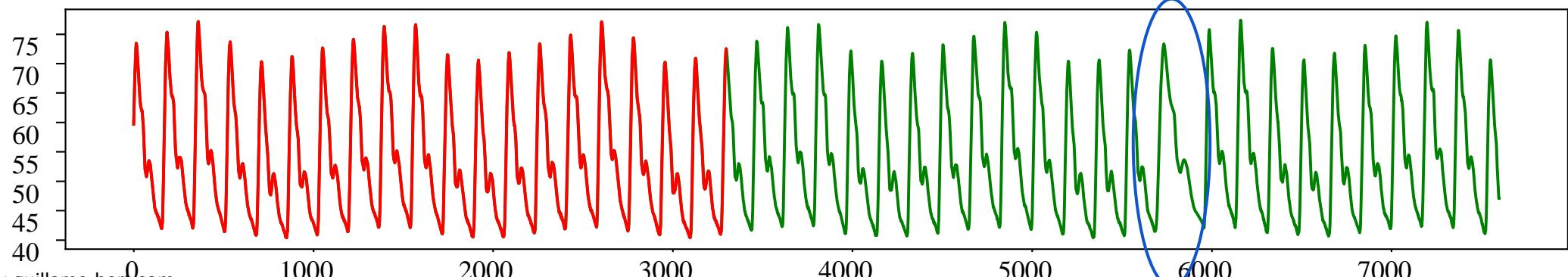
## UCR\_Anomaly\_InternalBleeding15\_1700\_5684\_5854.txt

The data comes from an internal bleeding dataset<sup>[1]</sup>. From the dataset, we extracted the arterial blood pressure measurements of pigs.

The anomaly is synthetic. We performed an upsampling operation on a short fraction of normal arterial blood pressure data. From 5684 to 5854, the frequency of time series is 1.6 times higher than before, and the missing data is filled with the value of the previous time point. This is a more subtle version of UCR\_Anomaly\_InternalBleeding6\_1500\_3474\_3629.txt.



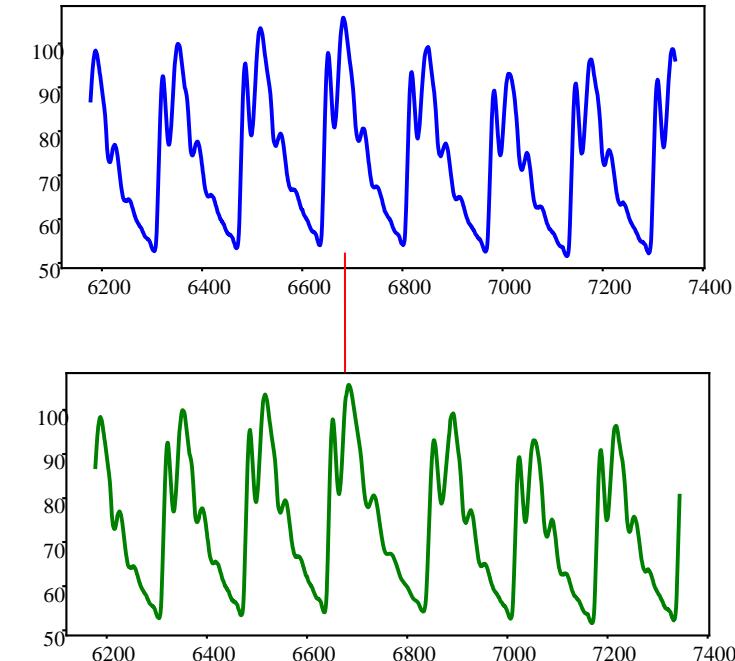
Blue is original data, green is data after anomaly was introduced



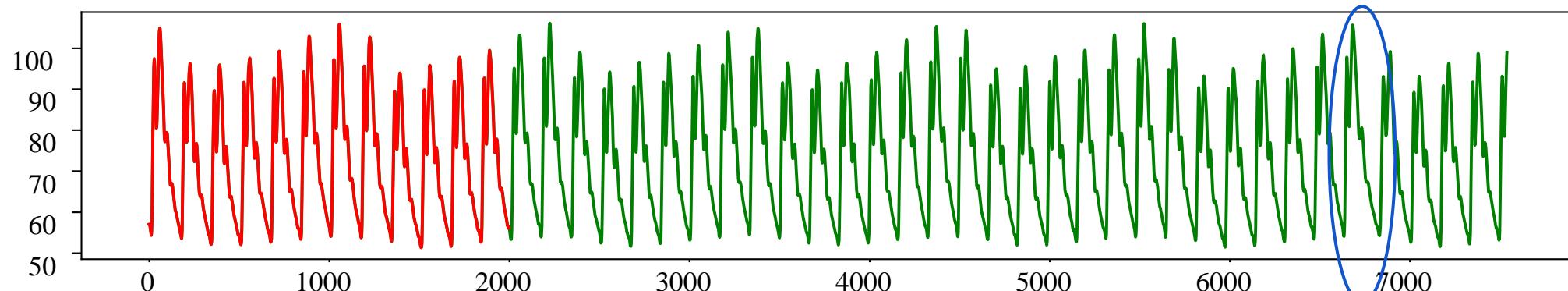
## UCR\_Anomaly\_InternalBleeding7\_2000\_6678\_6846.txt

The data comes from an internal bleeding dataset<sup>[1]</sup>. From the dataset, we extracted the arterial blood pressure measurements of pigs.

The anomaly is synthetic. We performed an upsampling operation on a short fraction of normal arterial blood pressure data. From 6678 to 6846, the frequency of time series is 1.25 times higher than before, and the missing data is filled with the value of the previous time point. This is a more subtle version of UCR\_Anomaly\_InternalBleeding6\_1500\_3474\_3629.txt.



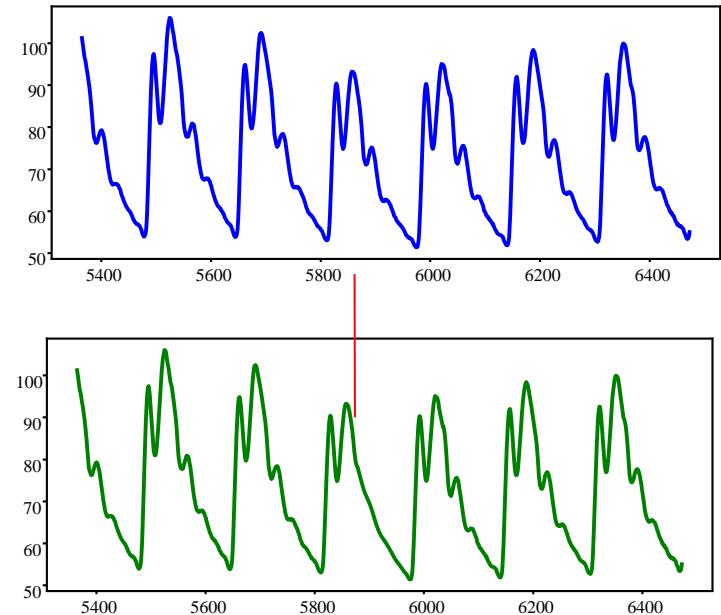
Blue is original data, green is data after anomaly was introduced



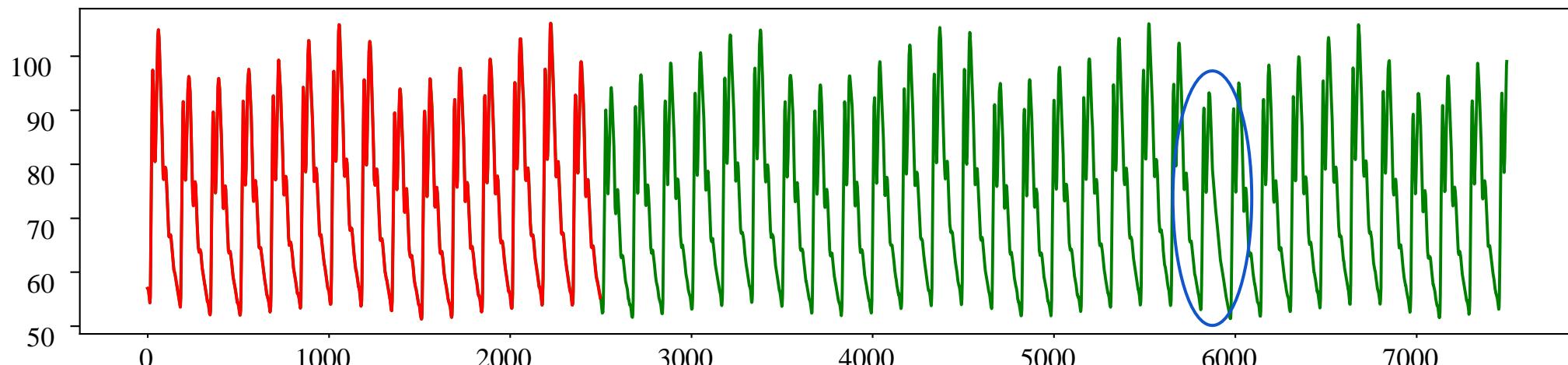
# UCR\_Anomaly\_InternalBleeding8\_2500\_5865\_5974.txt

The data comes from an internal bleeding dataset<sup>[1]</sup>. From the dataset, we extracted the arterial blood pressure measurements of pigs.

The anomaly is synthetic. We ran the moving average algorithm on a short fraction of normal arterial blood pressure data to eliminate about 4 extreme points within one cycle.



Blue is original data, green is data after  
anomaly was introduced

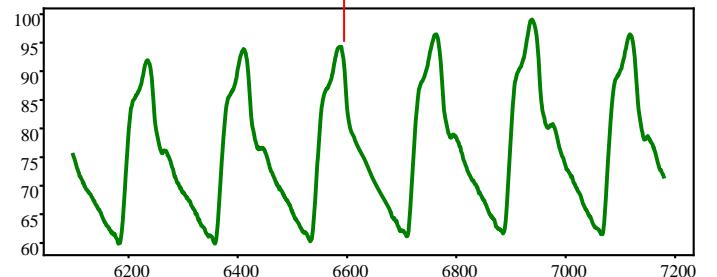
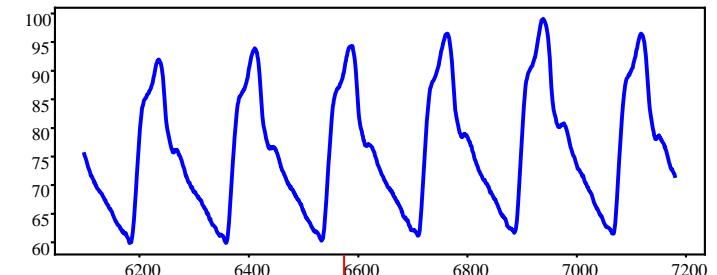


[1]: <http://mathieu.guillaume-bert.com>

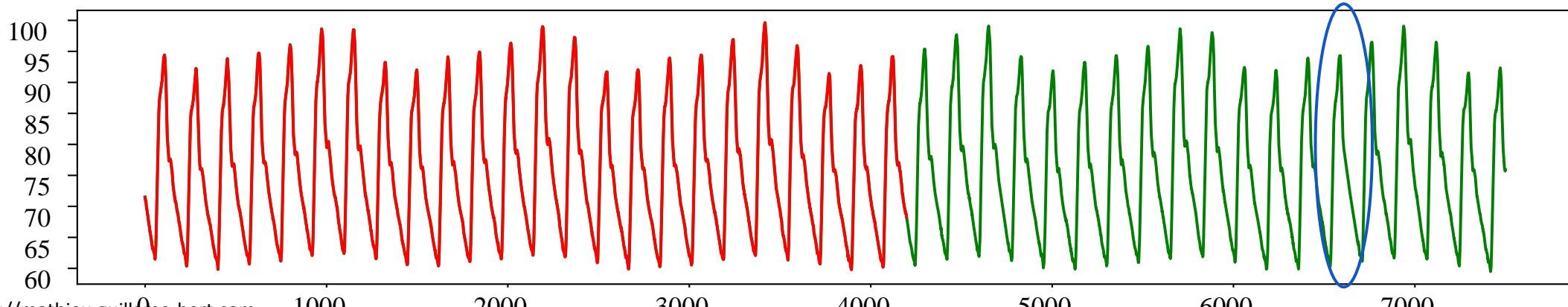
# UCR\_Anomaly\_InternalBleeding9\_4200\_6599\_6681.txt

The data comes from an internal bleeding dataset<sup>[1]</sup>. From the dataset, we extracted the arterial blood pressure measurements of pigs.

The anomaly is synthetic. We ran the moving average algorithm on a short fraction of normal arterial blood pressure data to eliminate about 2 extreme points within one cycle.



Blue is original data, green is data after anomaly was introduced

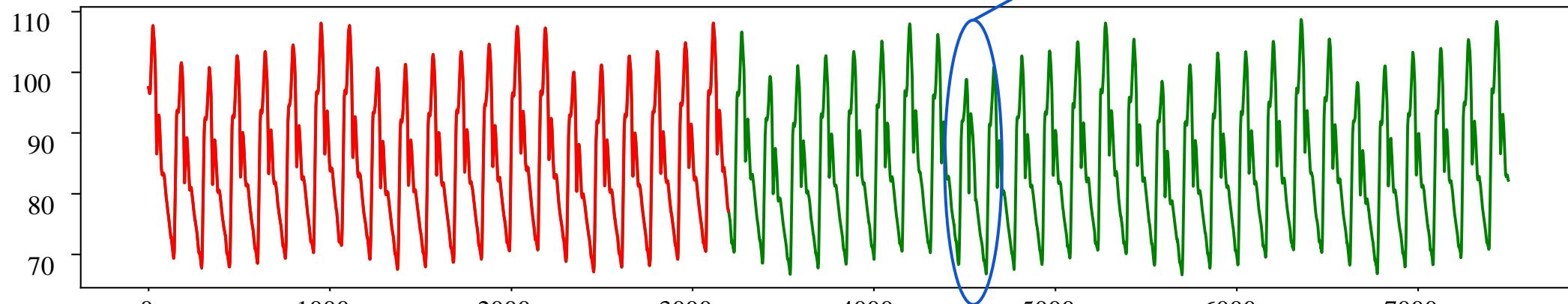
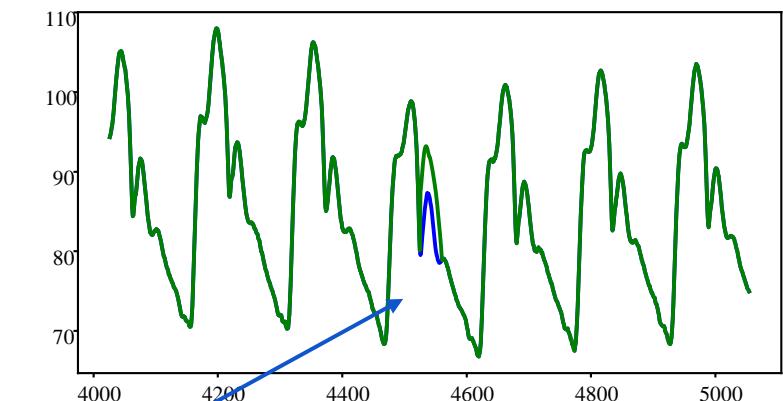


# UCR\_Anomaly\_InternalBleeding10\_3200\_4526\_4556.txt

The data comes from an internal bleeding dataset<sup>[1]</sup>. From the dataset, we extracted the arterial blood pressure measurements of pigs.

The anomaly is synthetic. We increased the amplitude of on a short fraction of normal arterial blood pressure data to 110% of the original to make the secondary peak in one cycle greater.

Blue is original data, green is data after anomaly was introduced

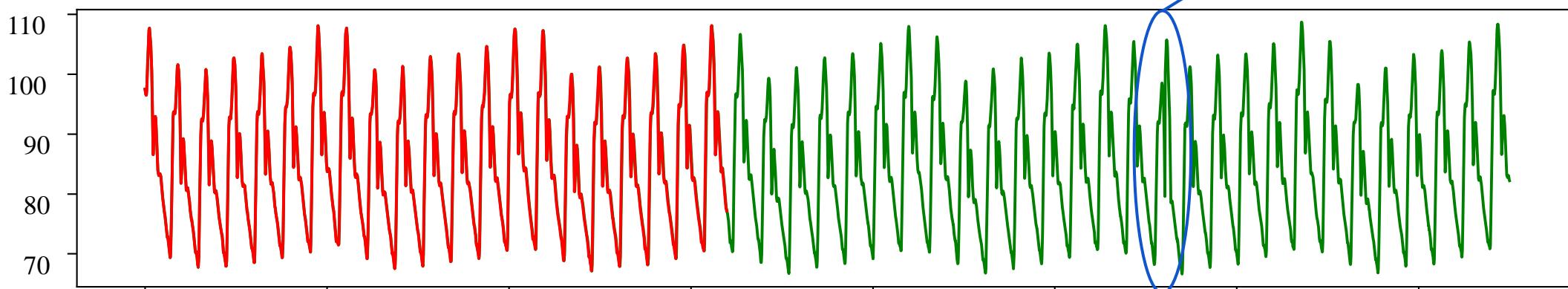
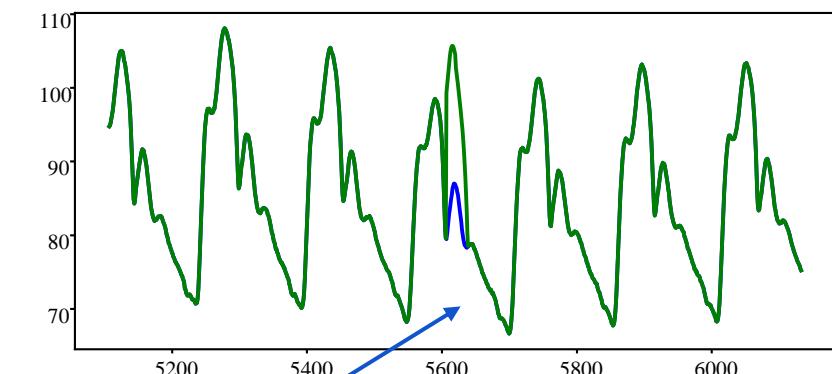


# UCR\_Anomaly\_InternalBleeding14\_2800\_5607\_5634.txt

The data comes from an internal bleeding dataset<sup>[1]</sup>. From the dataset, we extracted the arterial blood pressure measurements of pigs.

The anomaly is synthetic. We increased the amplitude of on a short fraction of normal arterial blood pressure data to **125%** of the original to make the secondary peak in one cycle greater than the main peak.

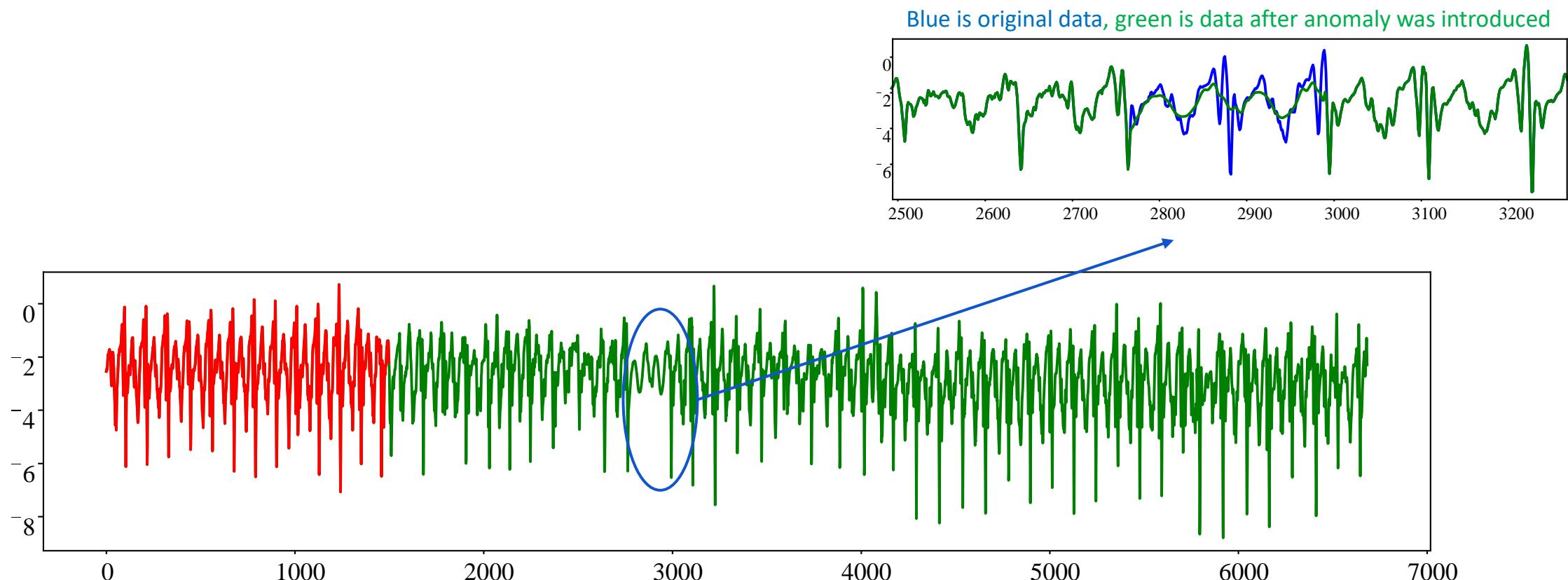
Blue is original data, green is data after anomaly was introduced



UCR Anomaly WalkingAcceleration1 1500 2764 2995.txt

The dataset is built based on an activity recognition system benchmark data<sup>[2]</sup>. We extracted an acceleration data measured by the sensor when an experimental subject is walking. Since the acceleration is 3-dimensional, we obtained our data by averaging the values in x, y and z axes.

The anomaly is synthetic. We first merged all data for the walking activity, and ran the moving average algorithm on a small fraction of the merged walking data (about two cycles).

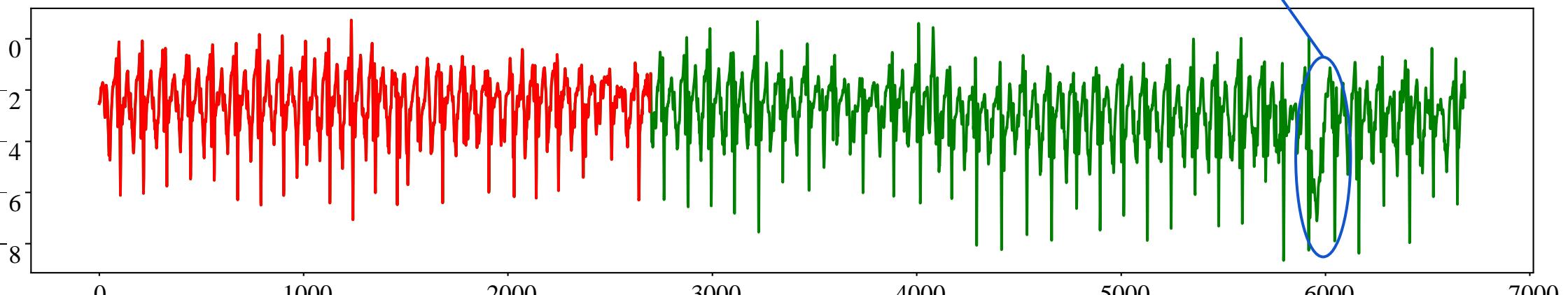
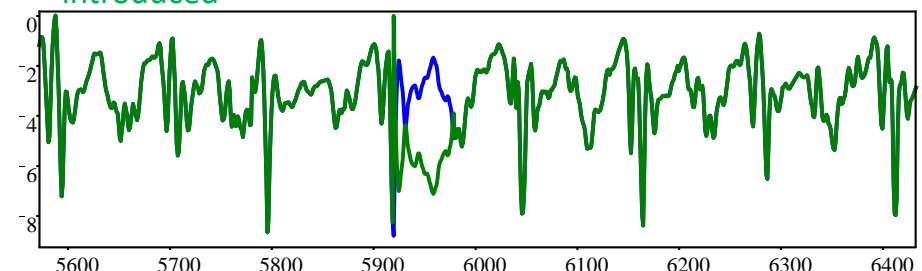


# UCR\_Anomaly\_WalkingAcceleration5\_2700\_5920\_5979.txt

The dataset is built based on an activity recognition system benchmark data<sup>[2]</sup>. We extracted an acceleration data measured by the sensor when an experimental subject is walking. Since the acceleration is 3-dimensional, we obtained our data by averaging the values in x, y and z axes.

We first merged all data for the walking activity. An anomaly is created by turning a half cycle's data upside down.

Blue is original data, green is data after anomaly was introduced



[2]: [https://www.dlr.de/kn/en/desktopdefault.aspx/tabcid-8500/14564\\_read-36508/](https://www.dlr.de/kn/en/desktopdefault.aspx/tabcid-8500/14564_read-36508/)

# UCR\_Anomaly\_GP711MarkerLFM5z1\_5000\_6168\_6212.txt

The data comes from subject 7 in the GaitPhase Database [1]. The subject was required to walk on a split-belt treadmill at walking speed of 1.1 m/s.

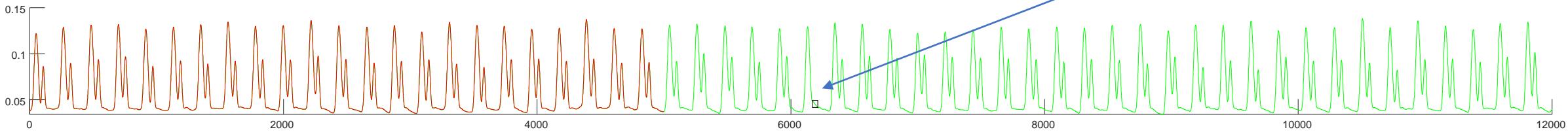
This dataset represents the vertical direction of a 3D marker's position. The marker was placed on subject's left shoe above the second metatarsal head.

There are totally 55 gait cycles in the dataset.

The anomaly is synthetic. For a random gait cycle, we created the anomaly by removing the second peak.

```
T(start_anomaly:end_anomaly) = interp1([start_anomaly end_anomaly],  
[T(start_anomaly) T(end_anomaly)], start_anomaly:end_anomaly);
```

Blue is original data, green is data after anomaly was introduced



# UCR\_Anomaly\_GP711MarkerLFM5z2\_5000\_5948\_5993.txt

The data comes from subject 7 in the GaitPhase Database [1]. The subject was required to walk on a split-belt treadmill at walking speed of 1.1 m/s.

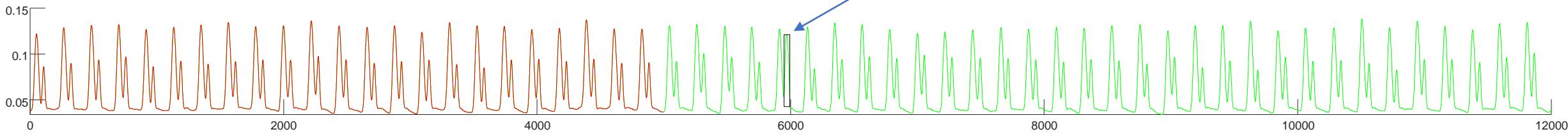
This dataset represents the vertical direction of a 3D marker's position. The marker was placed on subject's left shoe above the second metatarsal head.

There are totally 55 gait cycles in the dataset.

The anomaly is synthetic. For a random gait cycle, we made the second peak much higher, almost as high as the main peak.

```
[max_val,max_locidx] = max(T(start_anomaly:end_anomaly));
T(start_anomaly:end_anomaly) = interp1(
    [start_anomaly start_anomaly+max_locidx end_anomaly],
    [T(start_anomaly) max_val*1.25 T(end_anomaly)],
    start_anomaly:end_anomaly, 'pchip');
```

Blue is original data, green is data after anomaly was introduced



# UCR\_Anomaly\_GP711MarkerLFM5z3\_5000\_7175\_7388.txt

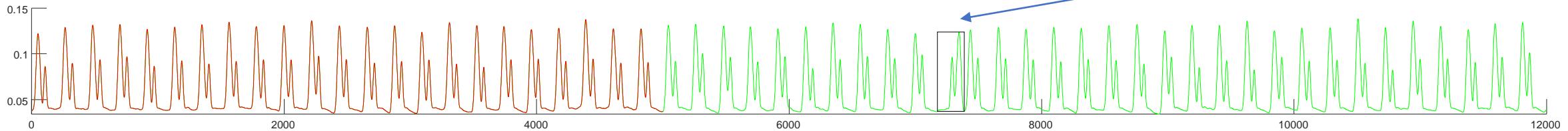
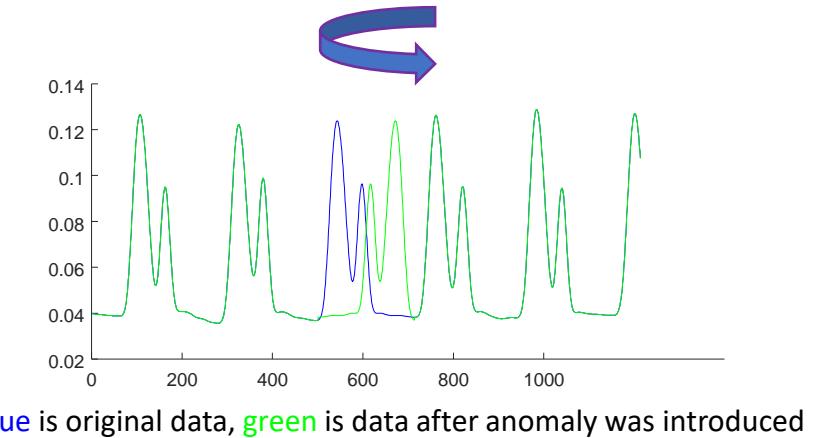
The data comes from subject 7 in the GaitPhase Database [1]. The subject was required to walk on a split-belt treadmill at walking speed of 1.1 m/s.

This dataset represents the vertical direction of a 3D marker's position. The marker was placed on subject's left shoe above the second metatarsal head.

There are totally 55 gait cycles in the dataset.

The anomaly is synthetic. We created the anomaly by reversing the direction of a random gait cycle.

```
T(start_anomaly:end_anomaly) = flip(T(start_anomaly:end_anomaly));
```



# UCR\_Anomaly\_GP711MarkerLFM5z4\_5000\_6527\_6645.txt

The data comes from subject 7 in the GaitPhase Database [1]. The subject was required to walk on a split-belt treadmill at walking speed of 1.1 m/s.

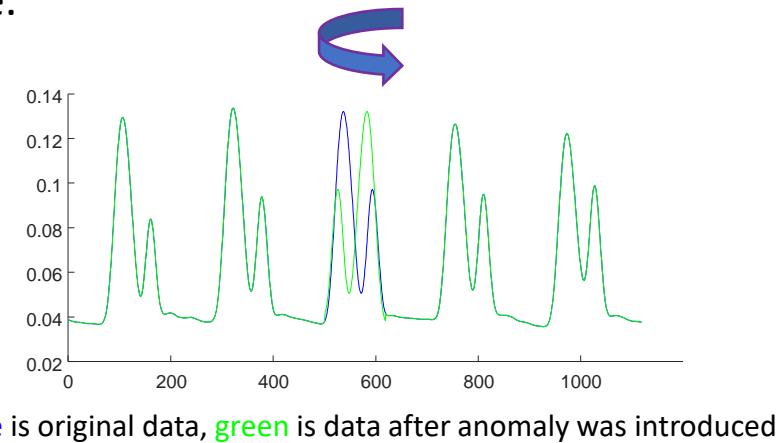
This dataset represents the vertical direction of a 3D marker's position. The marker was placed on subject's left shoe above the second metatarsal head.

There are totally 55 gait cycles in the dataset.

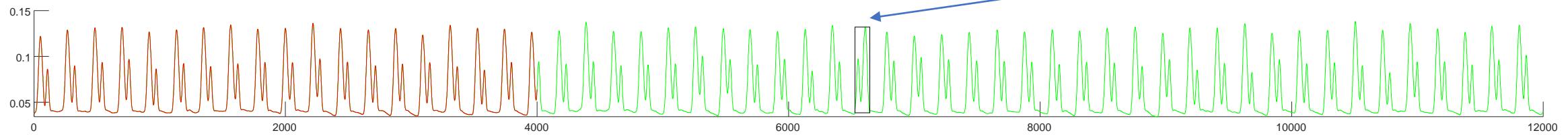
This is a subtle version of UCR\_Anomaly\_GP711MarkerLFM5z3\_5000\_7175\_7388.txt.

Instead of reversing the whole gait cycle, we only flipped the two peaks in the cycle.

```
T(start_anomaly:end_anomaly) = flip(T(start_anomaly:end_anomaly));
```



Blue is original data, green is data after anomaly was introduced



# UCR\_Anomaly\_GP711MarkerLFM5z5\_5000\_8612\_8716.txt

The data comes from subject 7 in the GaitPhase Database [1]. The subject was required to walk on a split-belt treadmill at walking speed of 1.1 m/s.

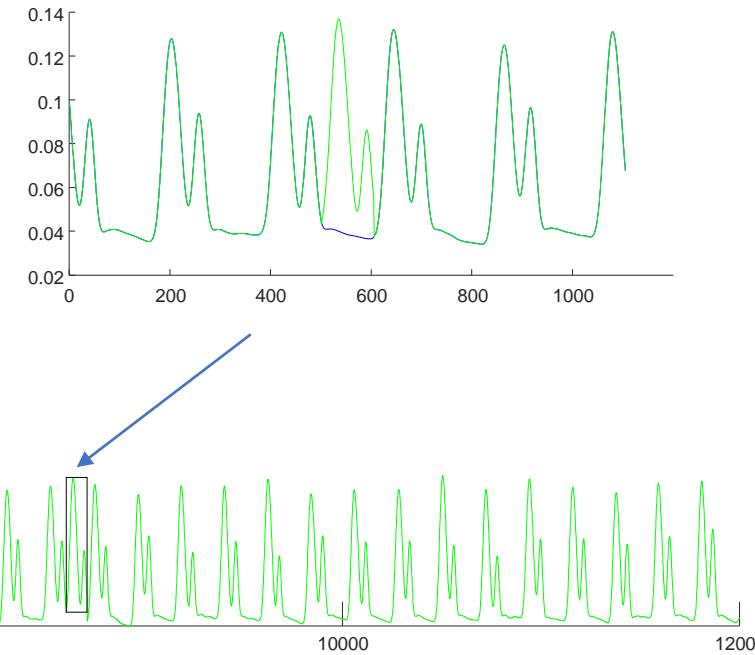
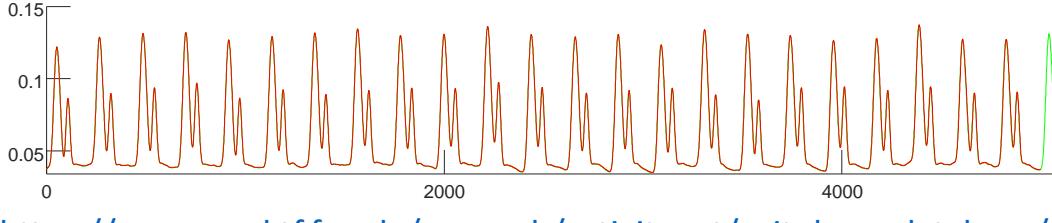
This dataset represents the vertical direction of a 3D marker's position. The marker was placed on subject's left shoe above the second metatarsal head.

There are totally 55 gait cycles in the dataset.

The anomaly is synthetic. For a random gait cycle, we replaced its swing phase with corresponding stance phase of subject's right foot.

Blue is original data, green is data after anomaly was introduced

```
T(start_anomaly:end_anomaly) = UCR_Original_GP711MarkerRFM5z(start_anomaly:end_anomaly);
```



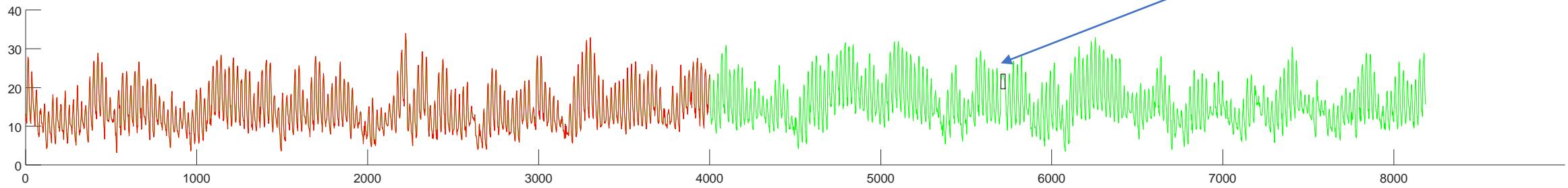
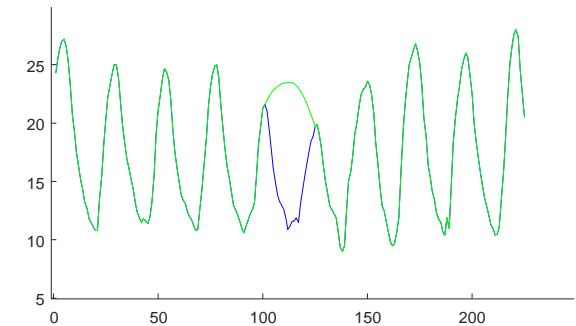
# UCR\_Anomaly\_CIMIS44AirTemperature2\_4000\_5703\_5727.txt

The data comes from the public weather data report from CIMIS station 44 in Riverside [1]. This dataset consists of hourly air temperature between 03/01 and 03/31 from 2009 to 2019.

The anomaly is synthetic. We removed the valley between randomly chosen 2 days and replaced the original two peaks with a single peak via cubic interpolation. The length of the original data is not changed but a wider peak was placed.

```
peak1 = T(start_anomaly);
peak2 = T(end_anomaly);
[~,valley_locidx] = min(T(start_anomaly:end_anomaly));
T(start_anomaly:end_anomaly) = interp1(
    [start_anomaly start_anomaly+valley_locidx end_anomaly],
    [peak1 3*max([peak1 peak2])-2*mean([peak1 peak2]) peak2],
    start_anomaly:end_anomaly, 'pchip');
```

Blue is original data, green is data after anomaly was introduced



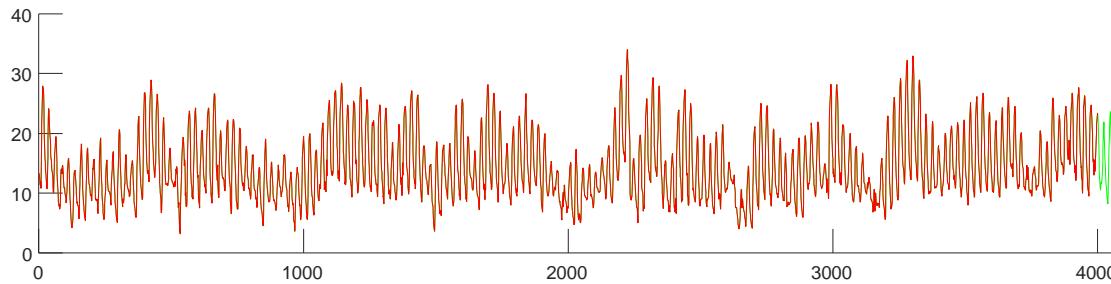
[1] <ftp://ftpcimis.water.ca.gov/pub2/annualMetric/>

# UCR\_Anomaly\_CIMIS44AirTemperature3\_4000\_6520\_6544.txt

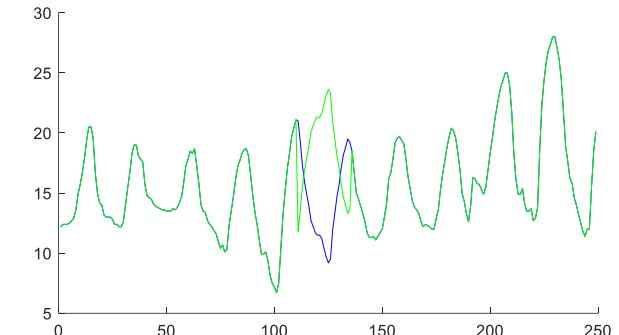
The data comes from the public weather data report from CIMIS station 44 in Riverside [1]. This dataset consists of hourly air temperature between 03/01 and 03/31 from 2009 to 2019.

The anomaly is synthetic. At a random location, we flipped a 12-hour data over across the mean value of the original data. The original valley now becomes an anomaly peak.

```
T(start_anomaly:end_anomaly) = 2*min(T(start_anomaly:end_anomaly))-  
T(start_anomaly:end_anomaly)+mean(T(start_anomaly:end_anomaly));
```



Blue is original data, green is data after anomaly was introduced



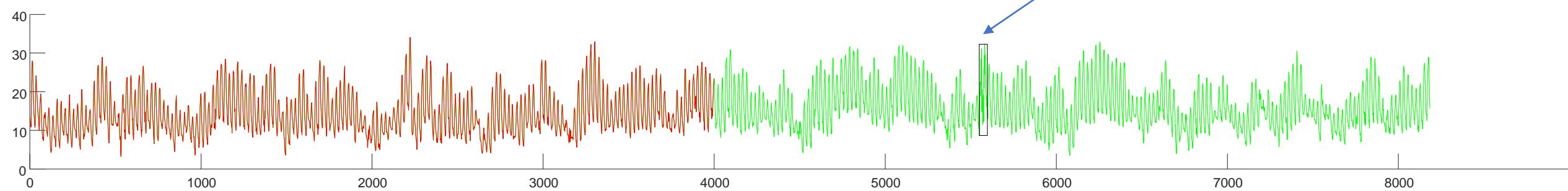
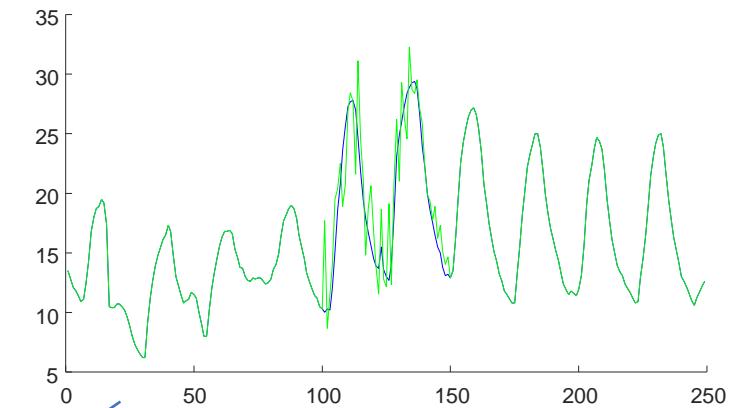
# UCR\_Anomaly\_CIMIS44AirTemperature4\_4000\_5549\_5669.txt

The data comes from the public weather data report from CIMIS station 44 in Riverside [1]. This dataset consists of hourly air temperature between 03/01 and 03/31 from 2009 to 2019.

The anomaly is synthetic. For a random consecutive 48 hours, we created the anomaly by adding noise, to simulate some electrical interference to the temperature sensor.

```
rng(223984);
T(start_anomaly:end_anomaly) = awgn(T(start_anomaly:end_anomaly),50,'measured','linear');
```

Blue is original data, green is data after anomaly was introduced



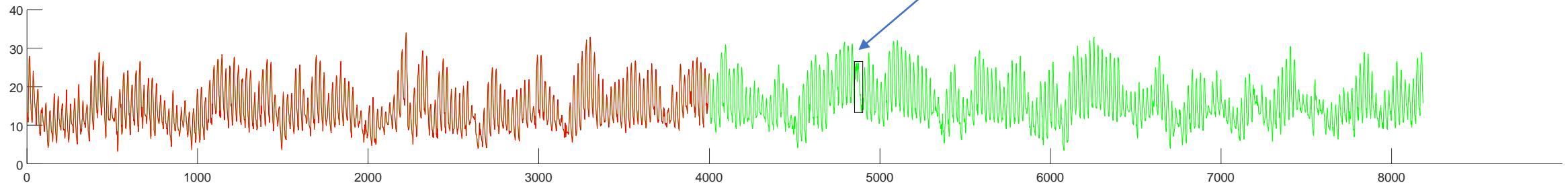
# UCR\_Anomaly\_CIMIS44AirTemperature5\_4000\_4852\_4900.txt

The data comes from the public weather data report from CIMIS station 44 in Riverside [1]. This dataset consists of hourly air temperature between 03/01 and 03/31 from 2009 to 2019.

The anomaly is synthetic. For a randomly chosen 48 hours, we replaced the original data with random walk to simulate a sudden failure of the temperature sensor. The random walk was carefully generated so that it fits within the value range of the original 48-hours data.

Blue is original data, green is data after anomaly was introduced

```
rng(485248);
min_temp = min(T(start_anomaly:end_anomaly));
max_temp = max(T(start_anomaly:end_anomaly));
anomaly = cumsum(randn(end_anomaly-start_anomaly+1, 1));
anomaly = (anomaly-min(anomaly))/(max(anomaly)-min(anomaly));
anomaly = anomaly*(max_temp-min_temp)+min_temp;
T(start_anomaly:end_anomaly) = anomaly;
```



## UCR\_Anomaly\_CIMIS44AirTemperature6\_4000\_6006\_6054.txt

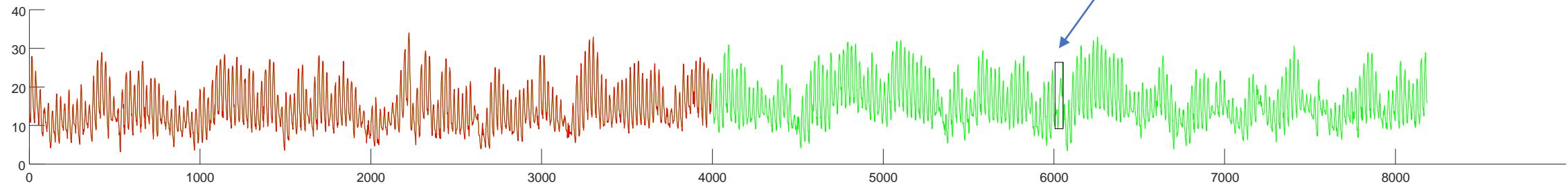
The data comes from the public weather data report from CIMIS station 44 in Riverside [1].

This dataset consists of hourly air temperature between 03/01 and 03/31 from 2009 to 2019.

The anomaly is synthetic. Same way as `UCR_Anomaly_CIMIS44AirTemperature4_4000_4852_4900.txt` to introduce the anomaly, but a much subtle version. Instead of directly replacing the original data with the random walk, we smoothed it before substitution to make it less obvious.

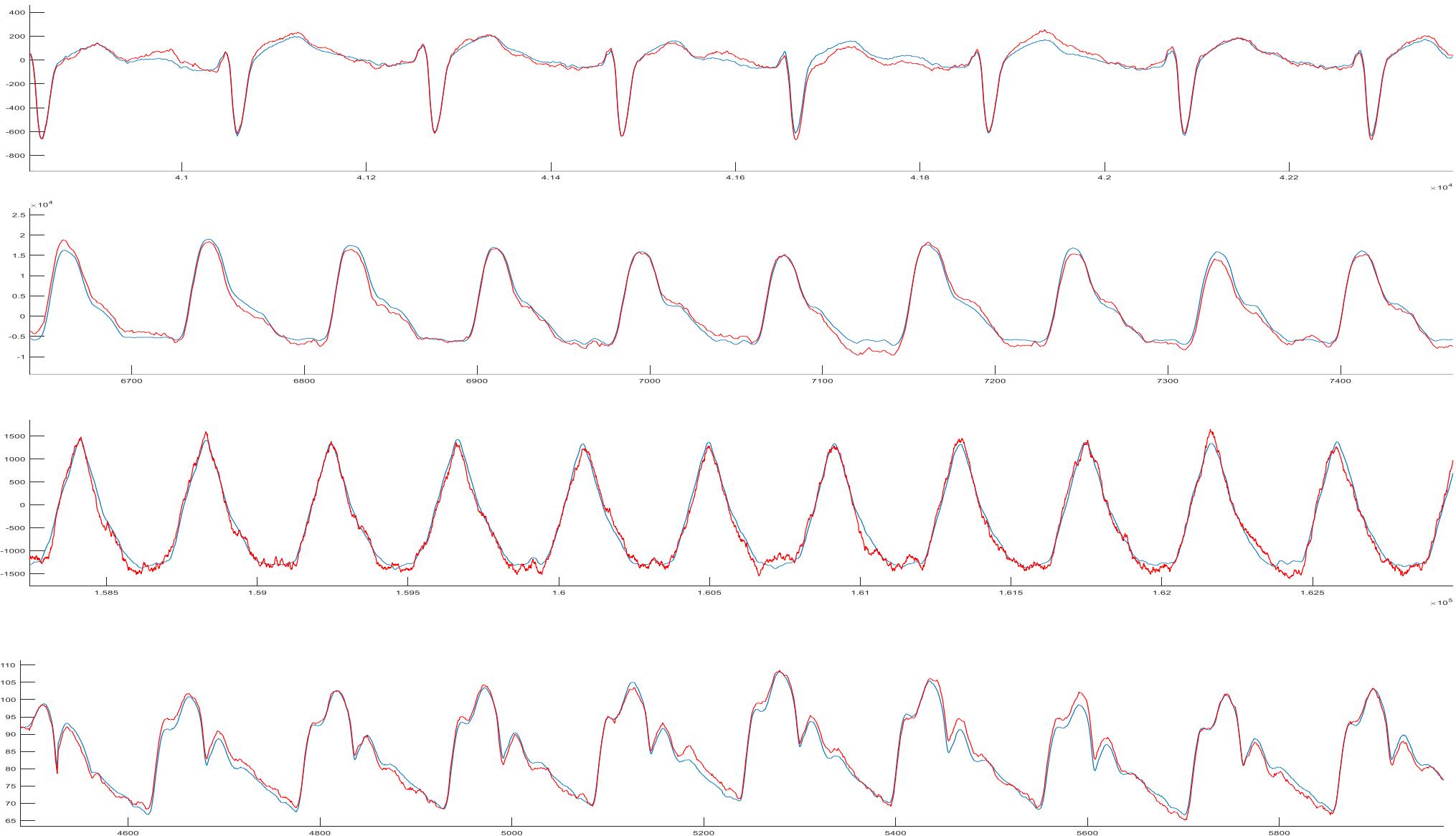
We also carefully chose the random seed, so that there is an abnormal peak temperature at dawn, which is quite impossible in real life. Because at dawn, it is usually the coolest time in a day.

```
rng(459728);
min_temp = min(T(start_anomaly:end_anomaly));
max_temp = max(T(start_anomaly:end_anomaly));
anomaly = cumsum(randn(end_anomaly-start_anomaly+1, 1));
anomaly = (anomaly-min(anomaly))/(max(anomaly)-min(anomaly));
anomaly = anomaly*(max_temp-min_temp)+min_temp;
T(start_anomaly:end_anomaly) = smooth(anomaly);
```

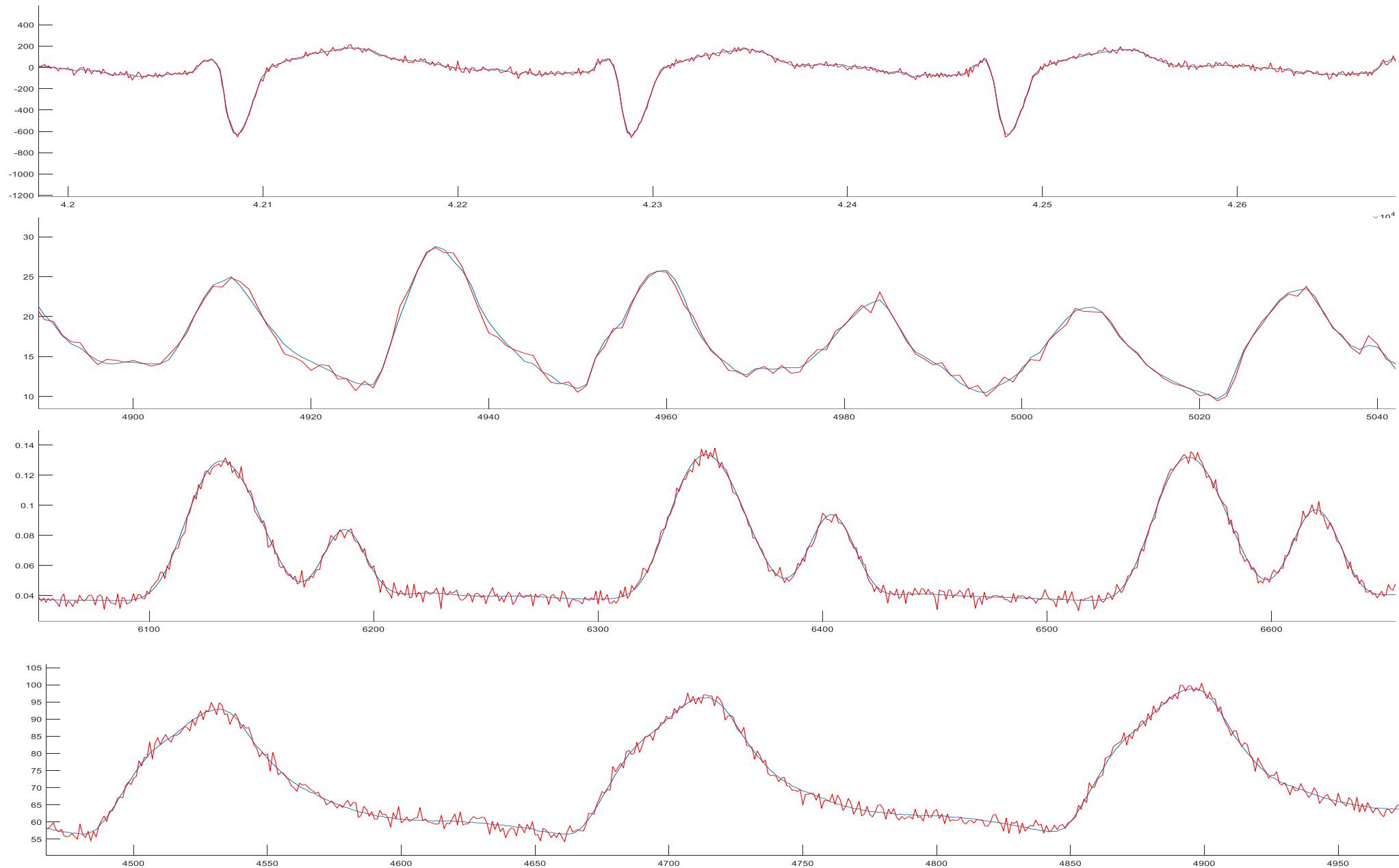


- The remaining datasets, upto 200, are duplicates of the previous datasets, but..
    - Some with some **distortion** added. If `temp` is the dataset, then...
      - `distort = temp + (smooth(randn(size(temp)), 25) * std(temp))`;
    - The name of the file reflects this. For example
      - `UCR_Anomaly_DISTORTEDECG4_5000_17000_17100.txt`
    - The distortion adds a little noise and wandering baseline.
  - and
  - Some with some **noise** added. If `temp` is the dataset, then...
    - `NOISE = temp + (randn(size(temp)) * (std(temp)/10))`;
  - The name of the file reflects this. For example
    - `102_UCR_Anomaly_NOISEMesoplodonDensirostris_10000_19280_19440.txt`
  - The noise is simple Gaussian noise
- In both cases, this makes little visual difference to the time series, see the next two slides.

- Samples of `distort = temp + (smooth(randn(size(temp)), 25) * std(temp))`; The red time series are distorted



- Samples of NOISE = temp + (randn(size(temp)) \* (std(temp)/10)); The red time series are noisy



For these files we considered data from MIT-BIH Long-term ECG database.

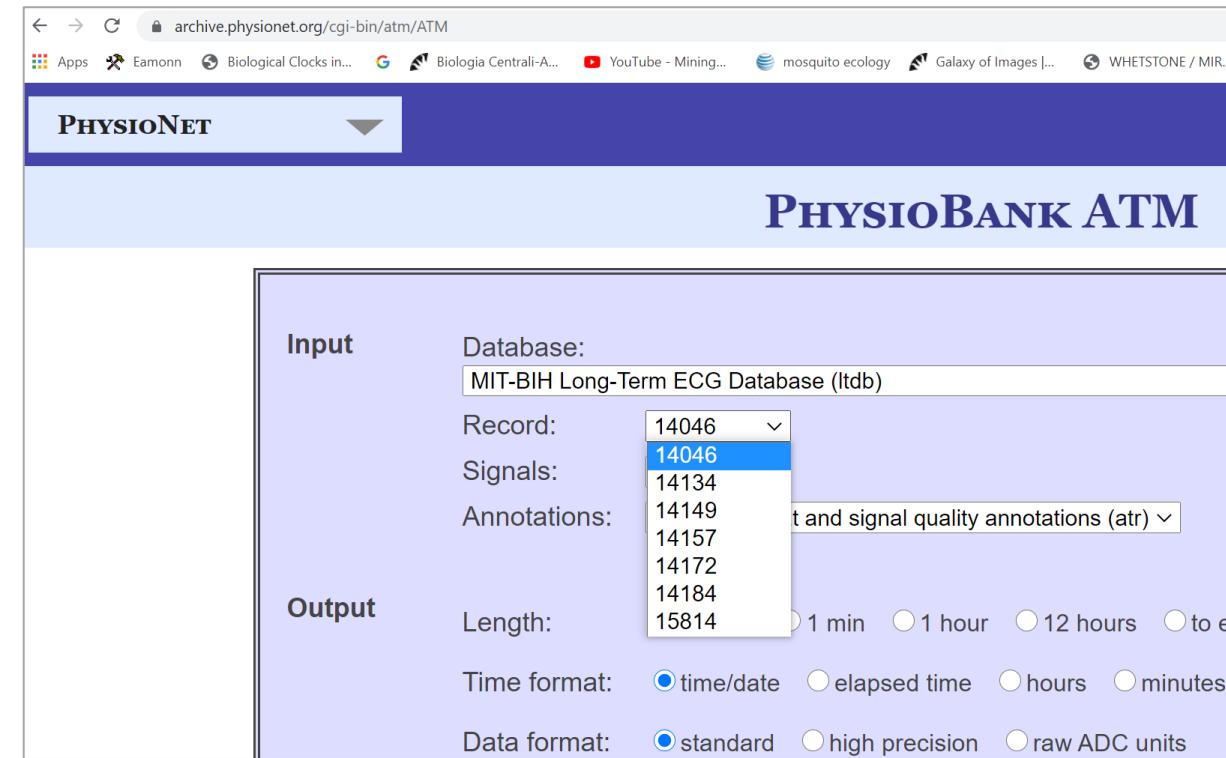
These dataset have many arrhythmias, sometimes multiple types of arrhythmias. However, in every case we made sure that the arrhythmias *also* exist in the training data, so they are not novel in the test data.

We will not bother creating a slide for each dataset, but in the following few slides we show some examples of the anomalies we created.

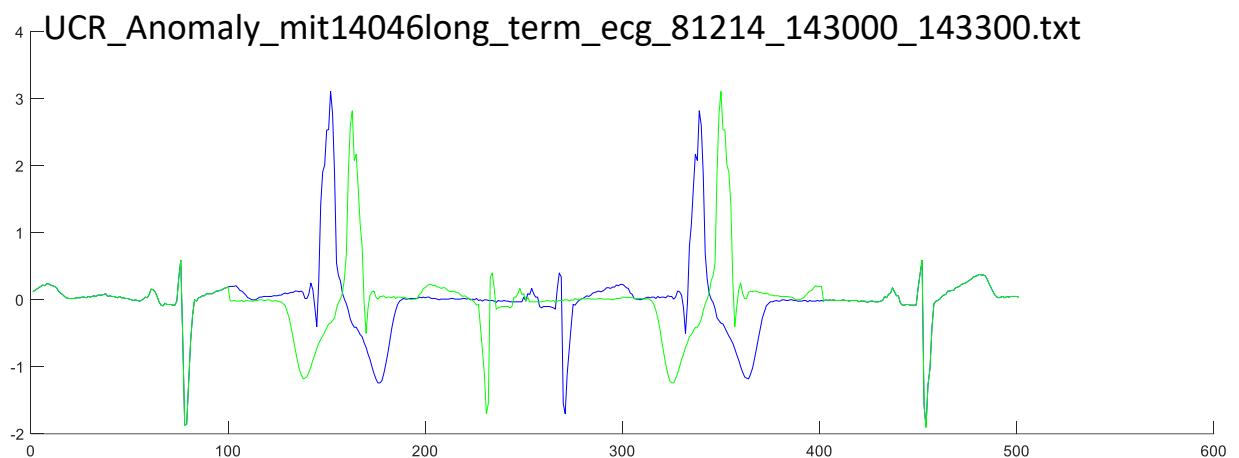
We create these with cardiologist Dr. Greg Mason

Our philosophy was: If a trained cardiologist had a lot of time to view the data, would *they* find the anomaly we added?

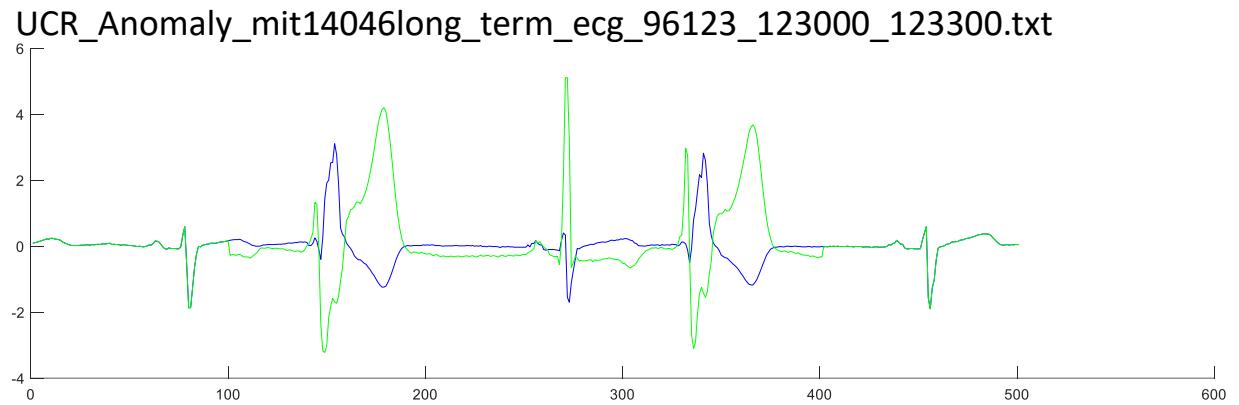
Some of the anomalies we added would easily pass this test, for some others it is closer to the edge...



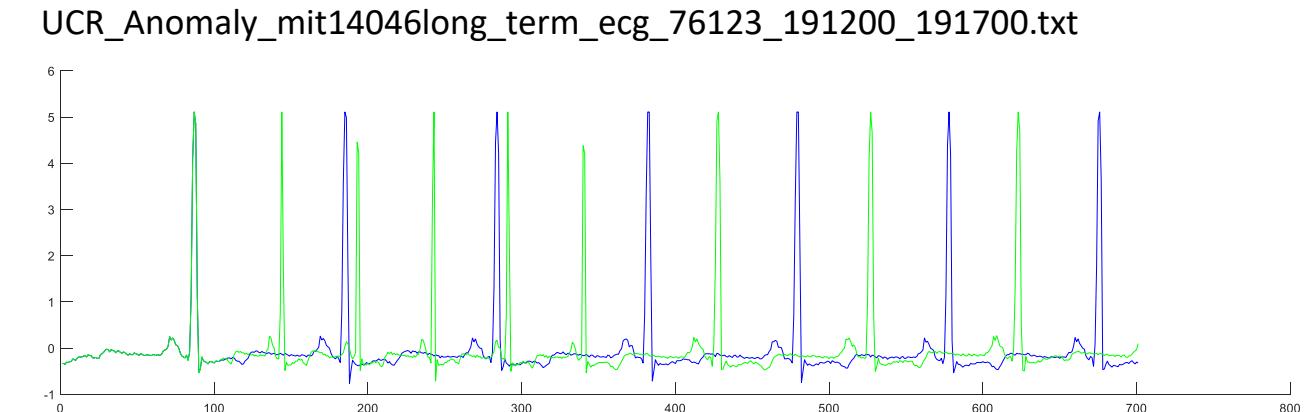
Reverse two beats (they happened to be arrythmias)



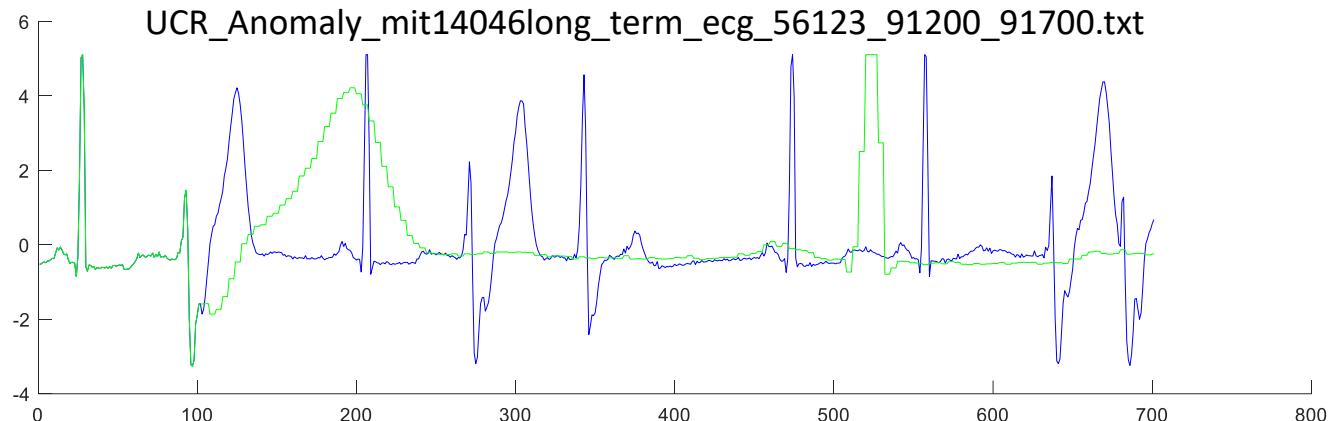
Swap in lead 2 for lead 1 a random region



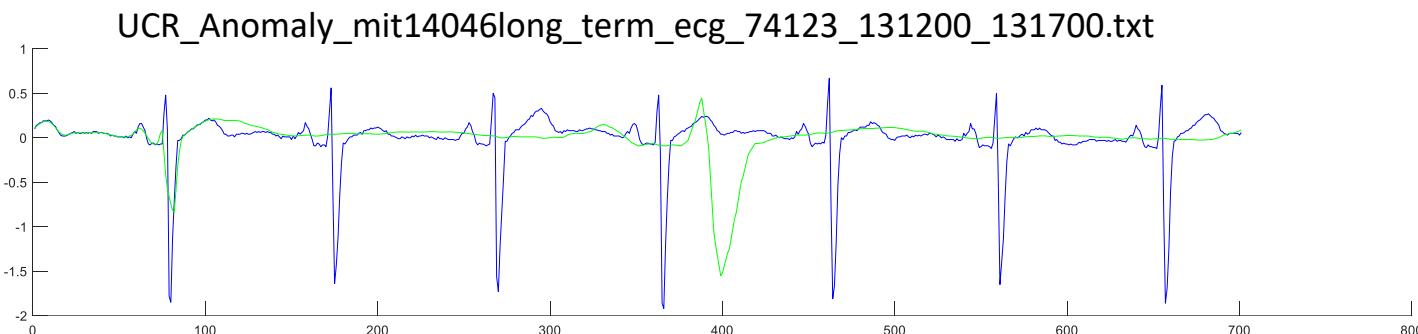
For a region containing about 8 normal beats, increase the speed by a factor of 2



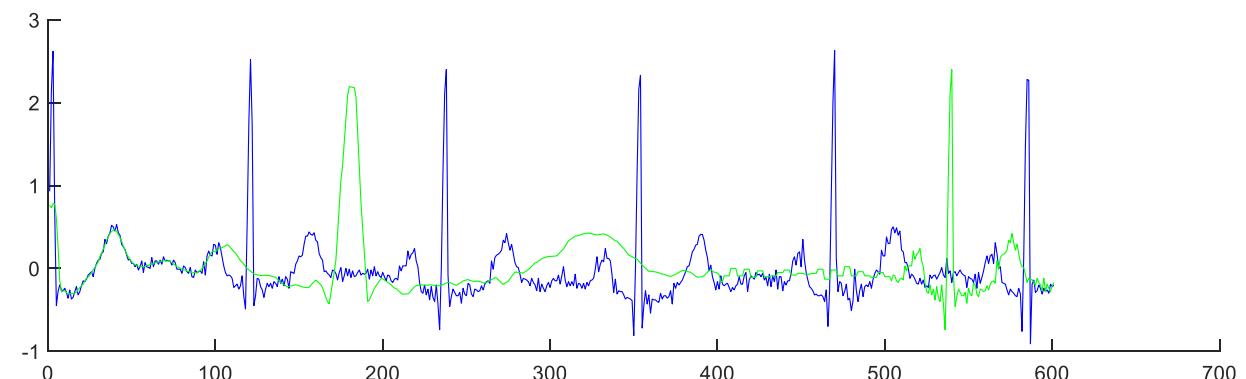
For a random region decrease the speed by a factor of 4, and leave in the “staircase” artifacts.



For a random region, decrease the speed by a factor of 4, and smooth out artifacts.

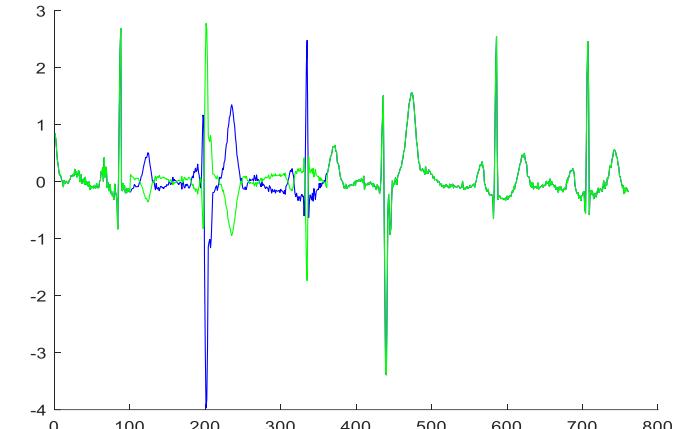


For a random region, decrease the speed by a factor of 4, and smooth out artifacts.



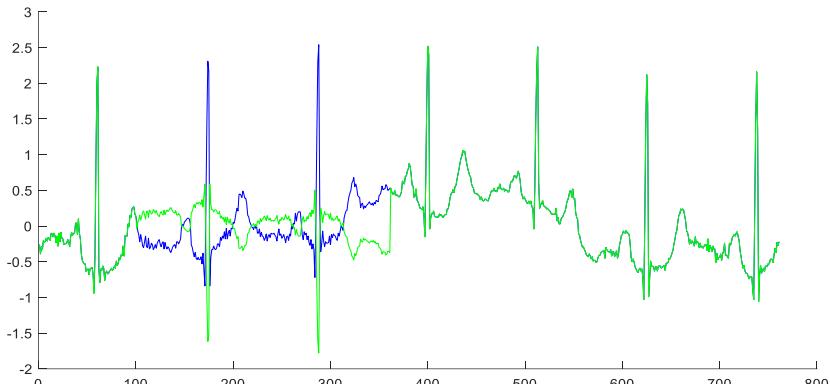
Flip a random region upside-down, by multiplying it by -0.7

UCR\_Anomaly\_mit14134long\_term\_ecg\_8763\_57530\_57790.txt



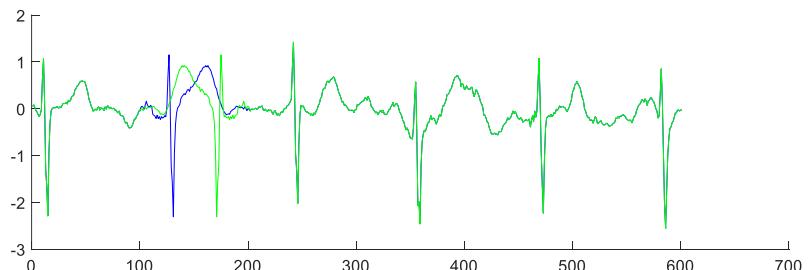
(again) Flip a random region upside-down, by multiplying it by -0.7

UCR\_Anomaly\_mit14134long\_term\_ecg\_8763\_47530\_47790.txt



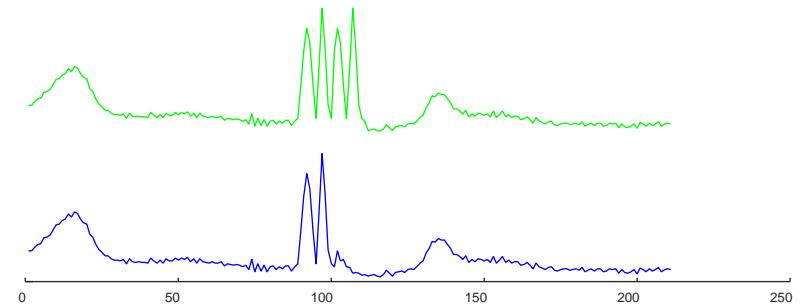
Reverse a single beat

UCR\_Anomaly\_mit14134long\_term\_ecg\_19363\_19510\_19610.txt



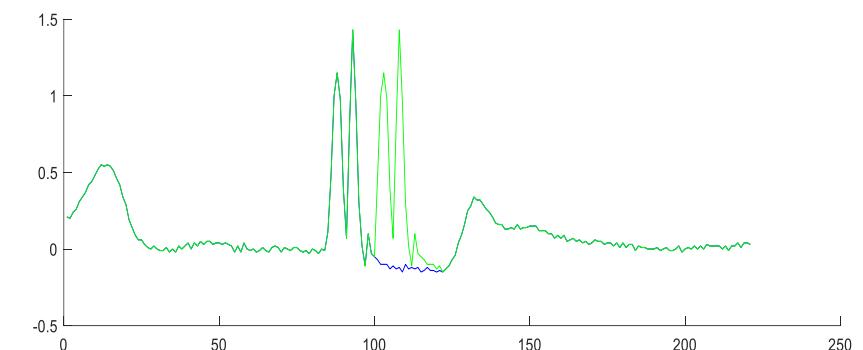
UCR\_Anomaly\_mit14134long\_term\_ecg\_16363\_57960\_57970.txt

Normally there are two peaks, we “stuttered” the data so that there are 4 peaks



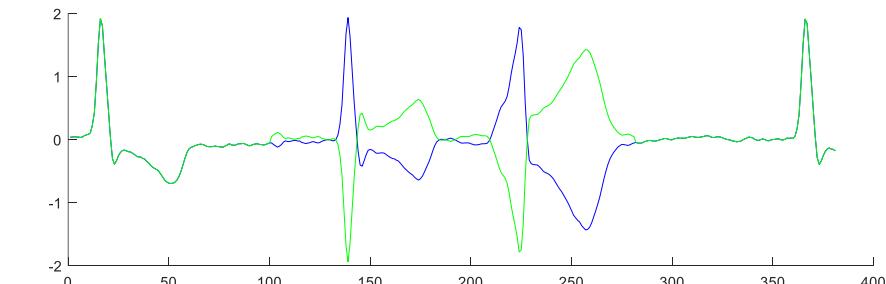
UCR\_Anomaly\_mit14134long\_term\_ecg\_11361\_47830\_47850.txt

(again) Normally there are two peaks, we “stuttered” the data so that there are 4 peaks

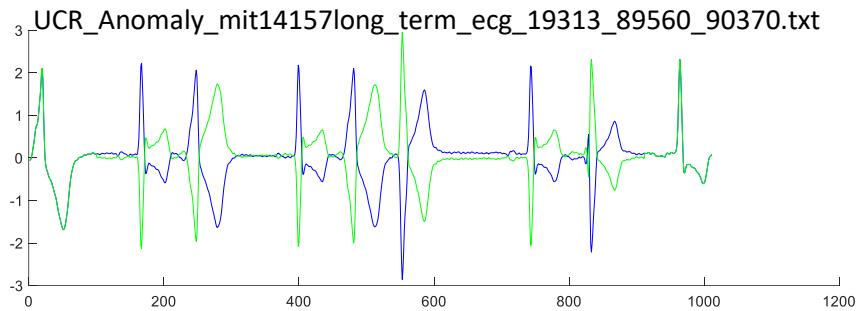


UCR\_Anomaly\_mit14157long\_term\_ecg\_21311\_72600\_72780.txt

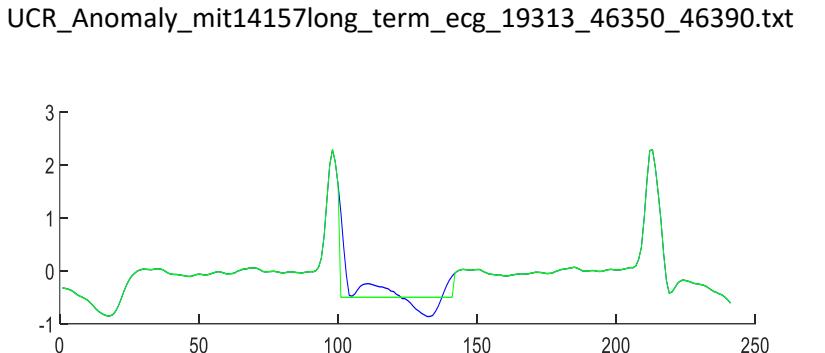
Flip a section upside down by multiplying it by -1



Flip a (longish) section that happened to contain some arrythmias upside down by multiplying it by -1



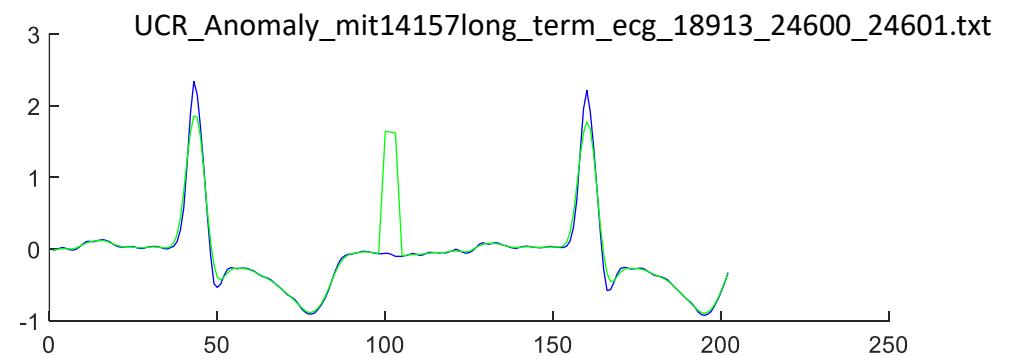
Hard limited a section of a beat to -0.5 (plus a tiny amount of noise)



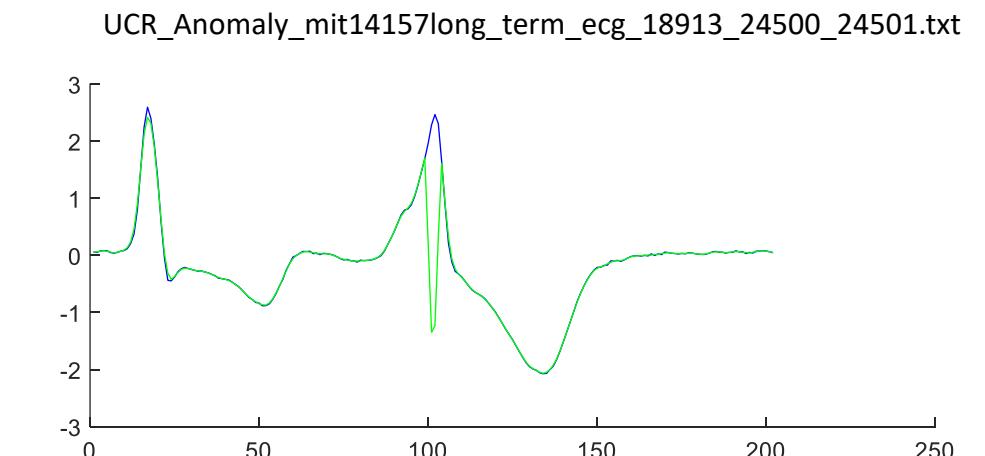
Added a spike (two points long)

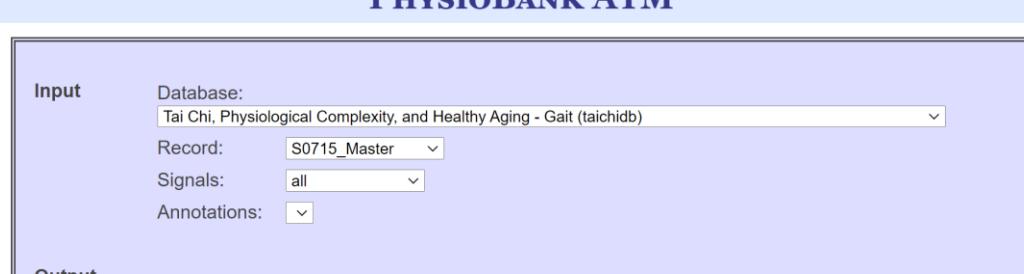


Added a spike (two points long), and slightly smoothed the area around it to blend it in more.



Added a dropout (two points long), and slightly smoothed the area around it to blend it in more.





The original data only has 4 unique values

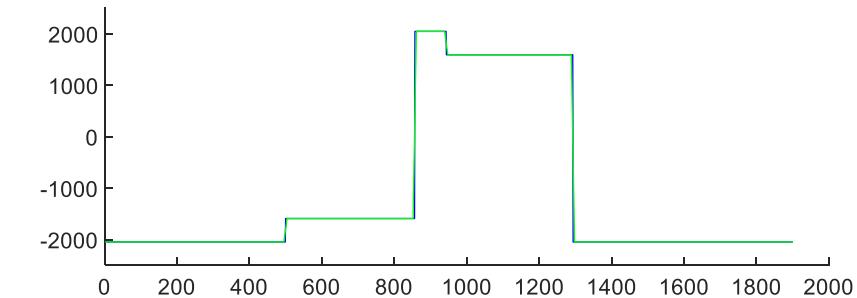
-2047    -1591    1592    2047

By *slightly* smoothing a section, we created a region with more unique values

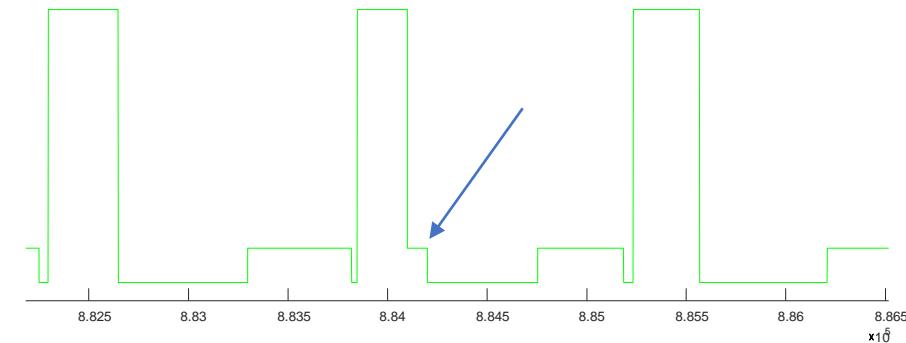
The original data has a small “step” before the main step. In one location we made a step after the main step.

The large step sometimes has an “overshoot” bump at the beginning. In one location we made an “overshoot” bump at the end of the step.

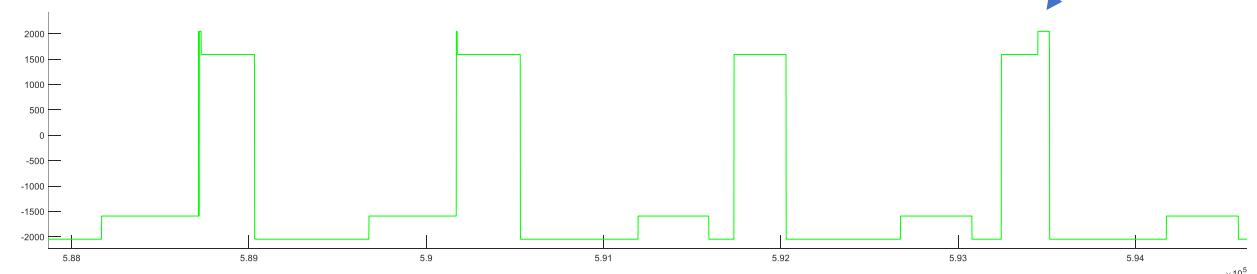
UCR\_Anomaly\_taichidbS0715Master\_250000\_837400\_839100.txt



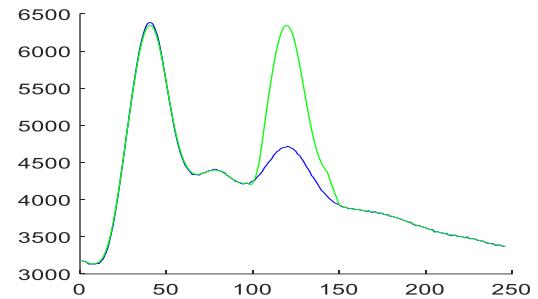
UCR\_Anomaly\_taichidbS0715Master\_240030\_884100\_884200.txt



UCR\_Anomaly\_taichidbS0715Master\_190037\_593450\_593514.txt

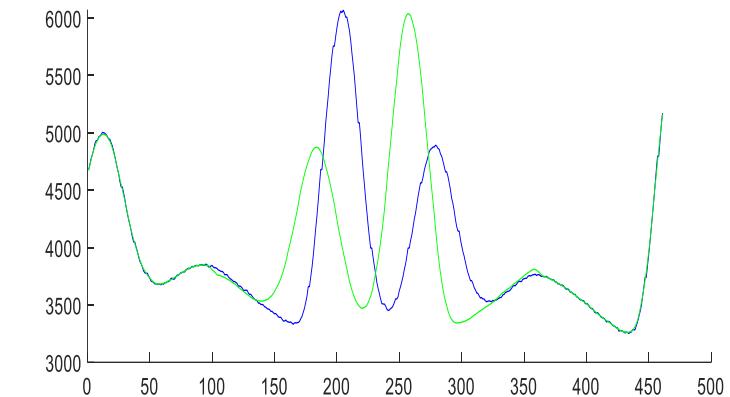


We “stuttered” the first bump of the APB cycle to enlarge the second bump



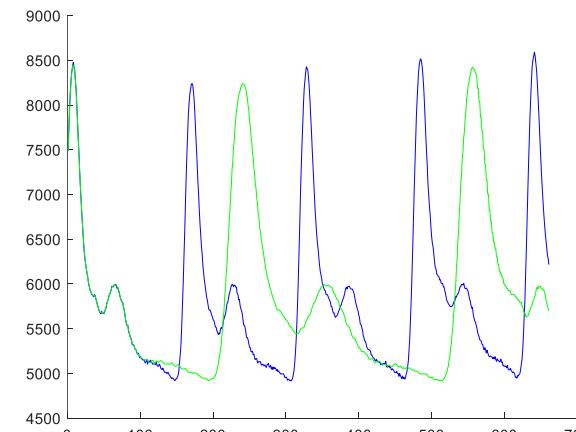
UCR\_Anomaly\_tilt12744mtable\_100000\_203355\_203400.txt

We reversed a single APB cycle



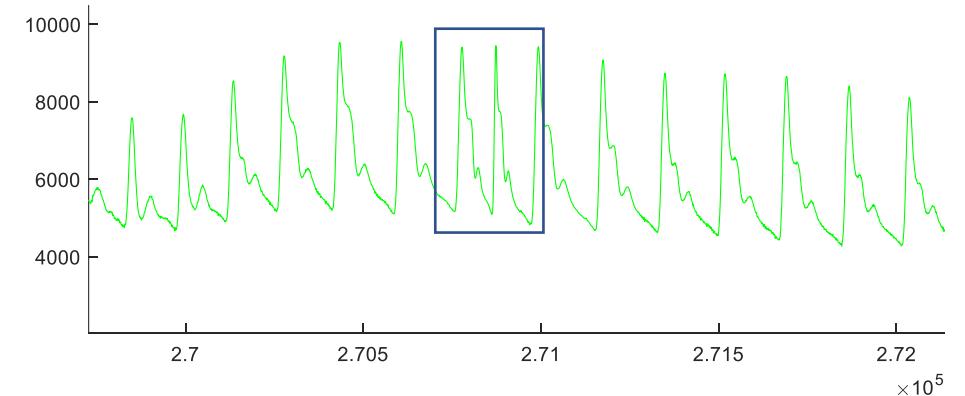
UCR\_Anomaly\_tilt12744mtable\_100000\_104630\_104890.txt

We slowed down an ABP beat by a factor of two.



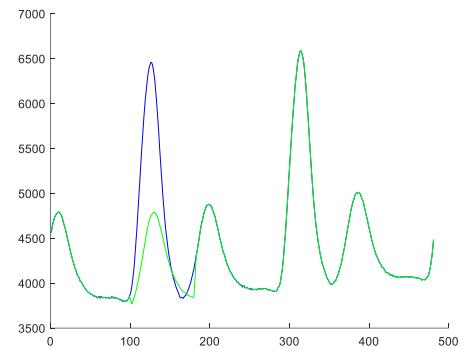
UCR\_Anomaly\_tilt12754table\_100013\_104630\_104890.txt

UCR\_Anomaly\_tilt12754table\_100211\_270800\_271070.txt



We sped up two ABP beats by a factor of two.

UCR\_Anomaly\_tilt12755mtable\_50211\_121900\_121980.txt



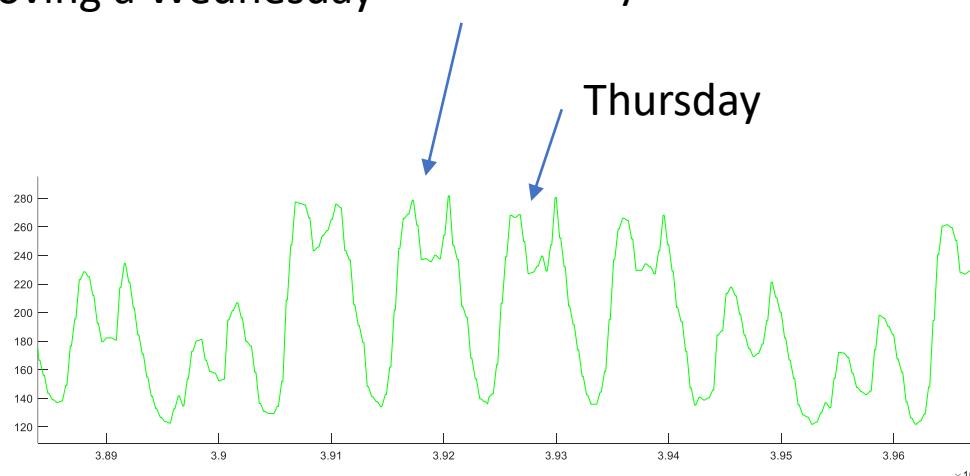
Stutter a small peak into the location of a large peak.

For Italian power demand

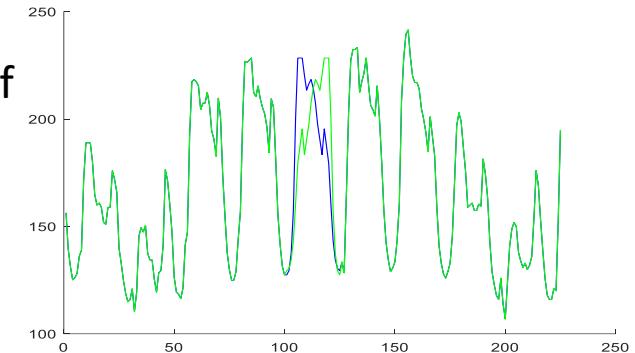
There are already natural “anomalies” (holidays, weather events etc). However, one year of training data should be enough to see all these events at least once.

So, the anomalies we embedded are obvious, and hopefully more significantly than any natural anomalies.

Made a 6-day week, by removing a Wednesday



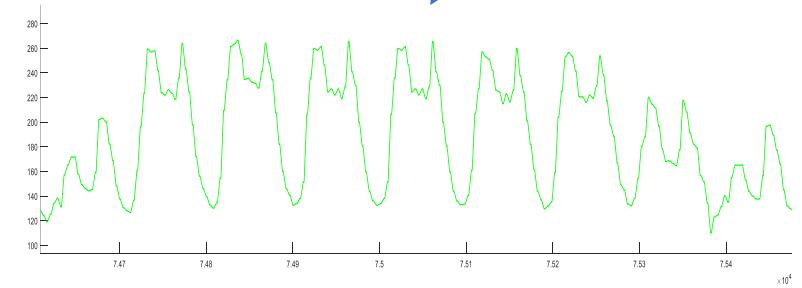
Reversed the direction of a Wednesday



UCR\_Anomaly\_Italianpowerdemand\_8913\_29480\_29504.txt

Made an 8-day week, by stuttering a Wednesday

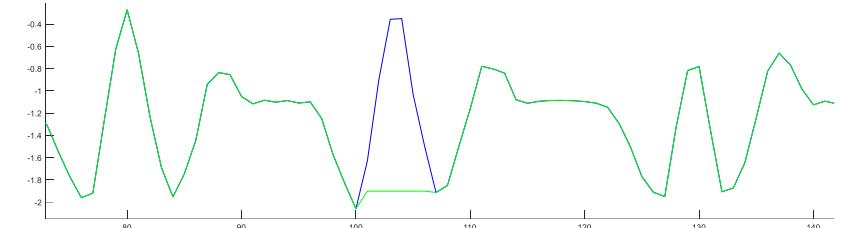
Real Wednesday  
Stuttered Wednesday



UCR\_Anomaly\_Italianpowerdemand\_36123\_74900\_74996.txt

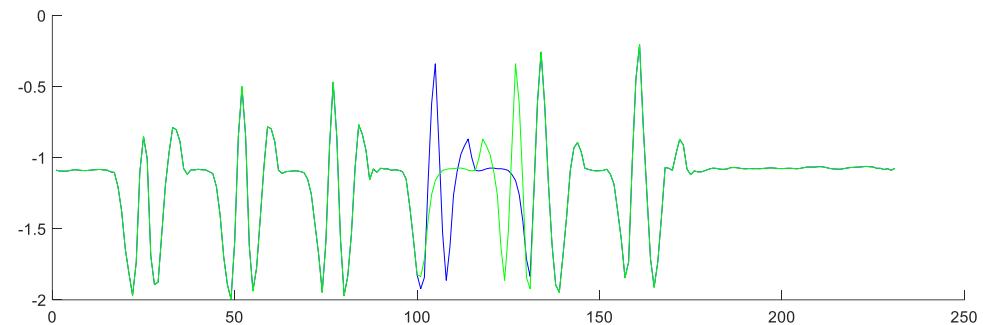
WeAllWalk

Flatted part  
of the gait  
cycle.



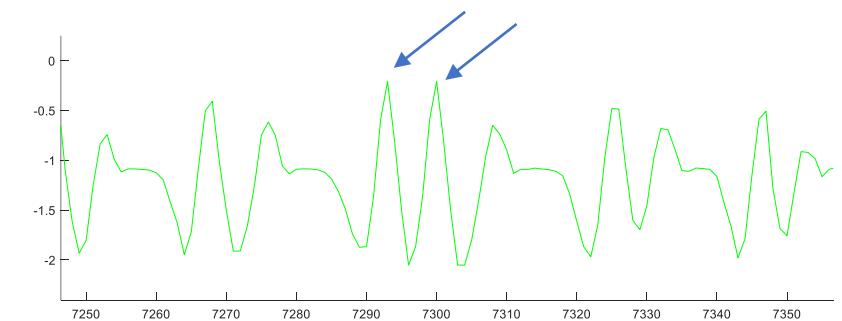
UCR\_Anomaly\_weallwalk\_2000\_4702\_4707.txt

Reverse 1  
gait cycle.



UCR\_Anomaly\_weallwalk\_2753\_8285\_8315.txt

“stuttered”  
part of gait  
cycle



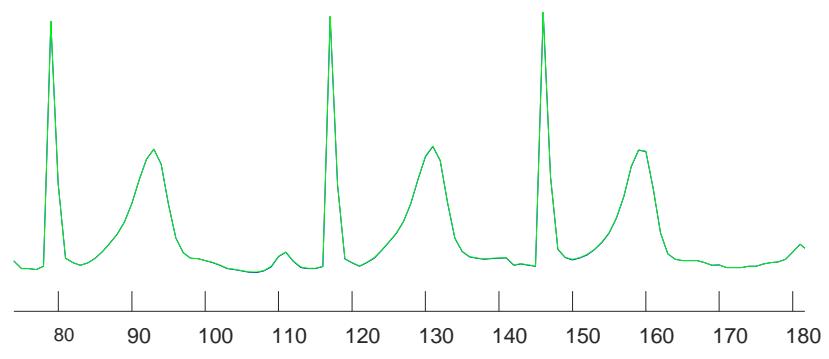
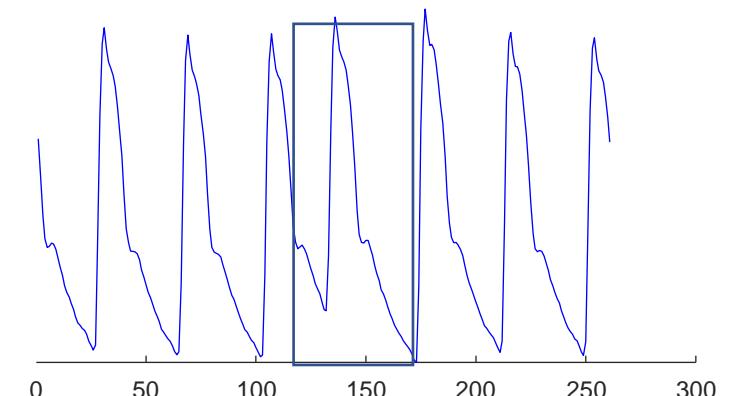
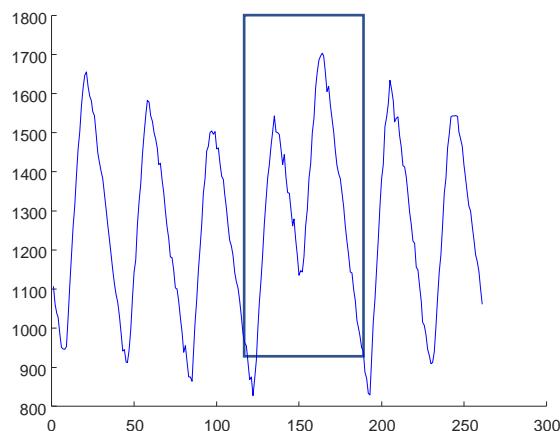
UCR\_Anomaly\_weallwalk\_2951\_7290\_7296.txt

## CHARISten

ECG, arterial blood pressure (ABP), and intracranial pressure (ICP)

Natural anomaly of a single arrhythmia.

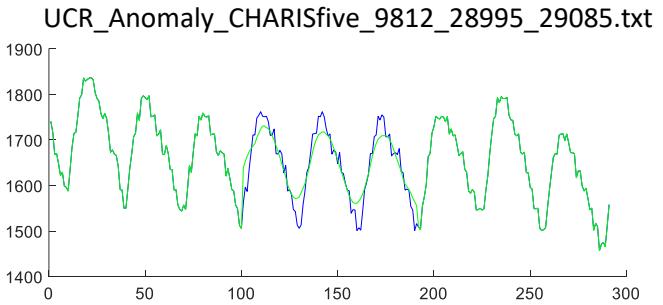
We shifted the starting point, so that the anomaly would not always be in the same spot for the 3 datasets we made.



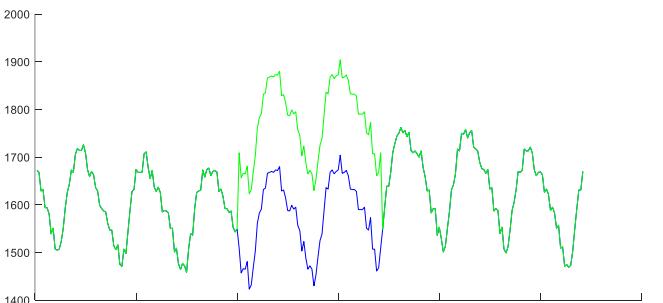
## CHARISfive

ECG, arterial blood pressure (ABP), and intracranial pressure (ICP)

We smoothed three beats, which also has the effect of reducing the size of the peaks and valleys



We added 200 to two beats, shifting them up.



This dataset

202\_UCR\_Anomaly\_CHARISfive\_10411\_10998\_11028.txt

Was originally mislabeled as

202\_UCR\_Anomaly\_CHARISfive\_11411\_10998\_11028.txt

Meaning that it was impossible to solve.

It is now fixed.



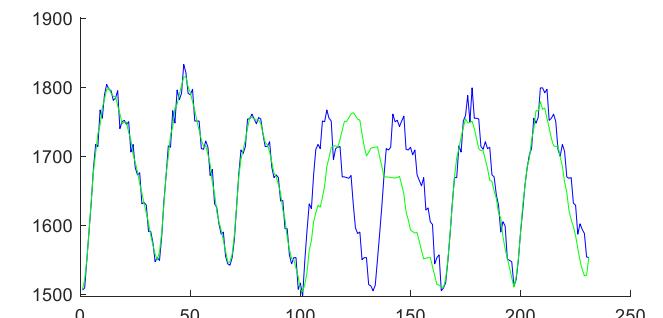
We made one beat twice as long as normal

Same for

203\_UCR\_Anomaly\_CHARISfive\_11812\_10995\_11028.txt

Now fixed, and renamed to

203\_UCR\_Anomaly\_CHARISfive\_10500\_10995\_11028.txt

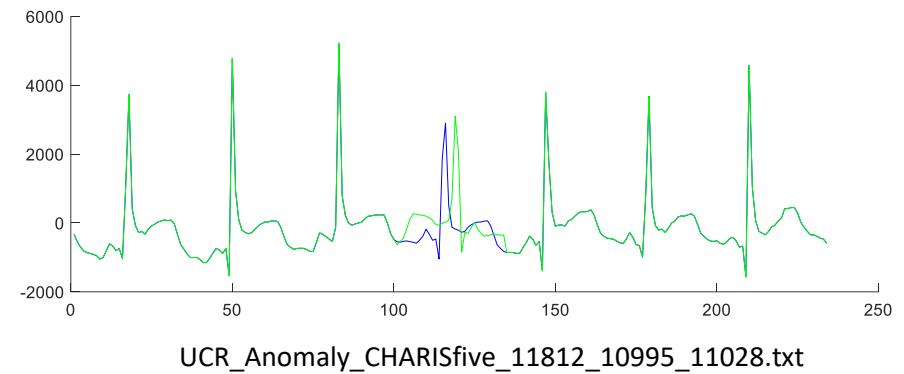


UCR\_Anomaly\_CHARISfive\_11411\_10998\_11028.txt

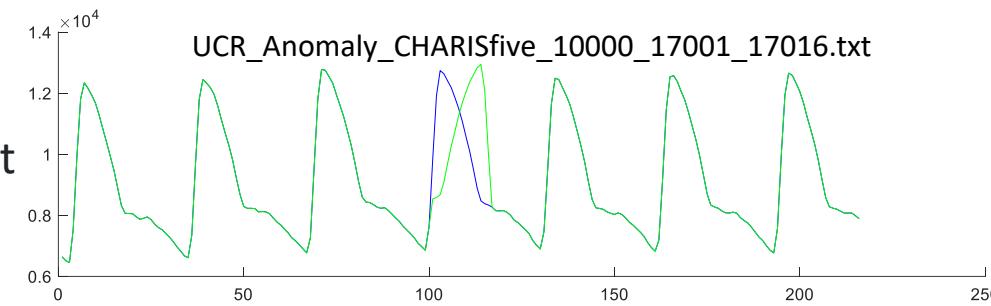
CHARISfive

ECG, arterial blood pressure (ABP), and  
intracranial pressure (ICP)

We flipped a beat

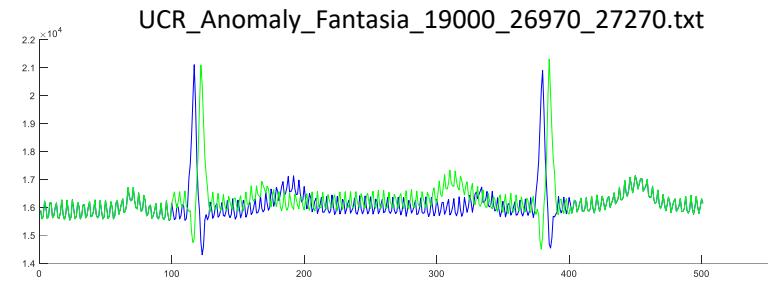


We flipped a beat



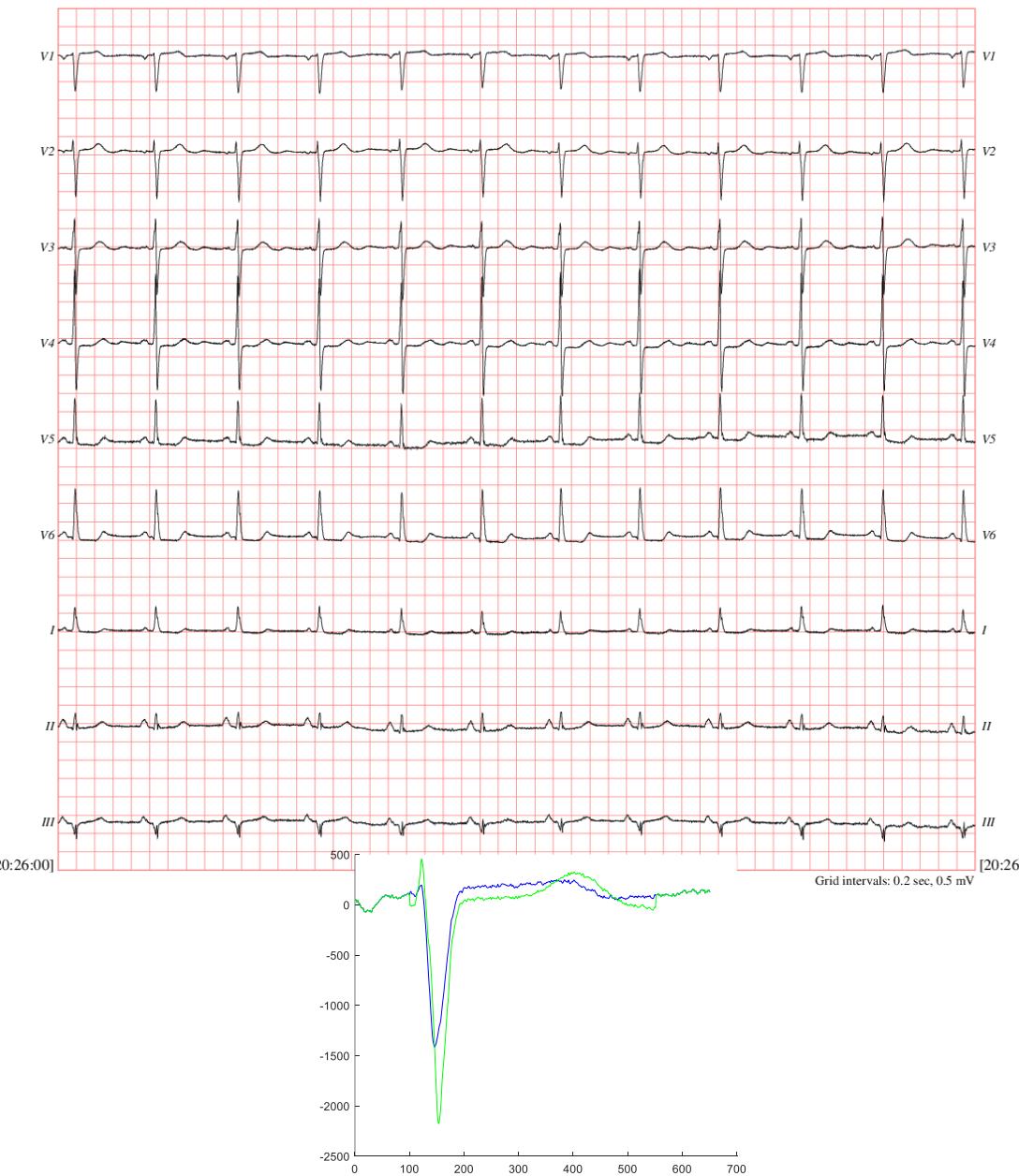
# Fantasia Database

We flipped a beat



# STAFFIII Database

Here we swapped in short regions from a different trace



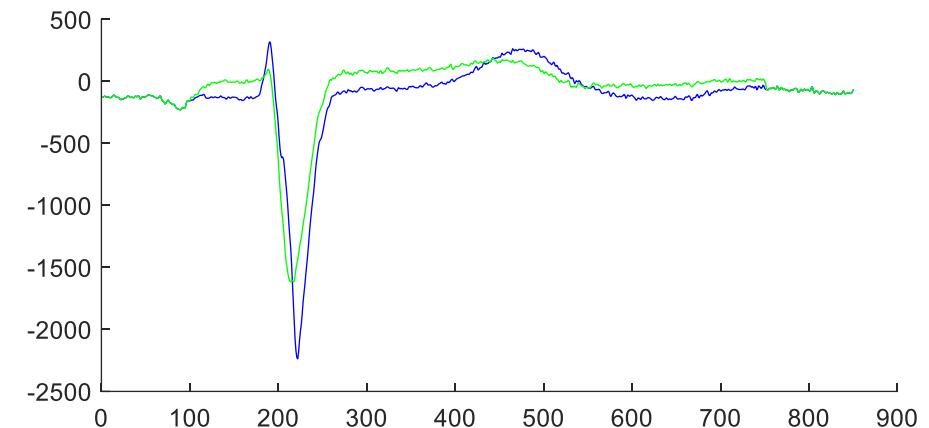
V2 for V1

UCR\_Anomaly\_STAFFIII Database\_33211\_126920\_127370.txt

## STAFFIIDatabase

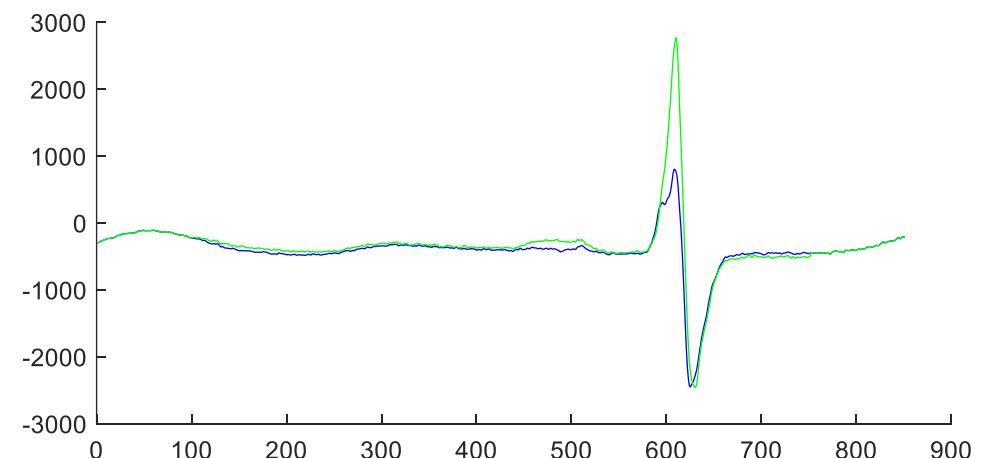
Here we swapped  
in short regions  
from a different  
trace

V2 for V1



UCR\_Anomaly\_STAFFIIDatabase\_36276\_106720\_107370.txt

V3 for 4

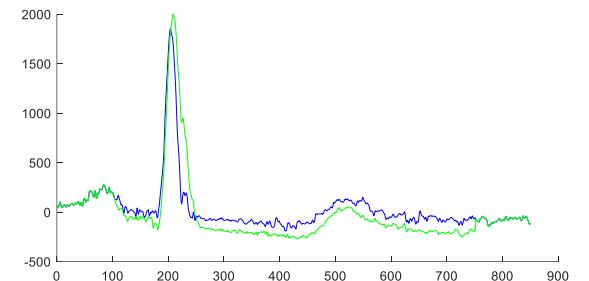


UCR\_Anomaly\_STAFFIIDatabase\_37216\_160720\_161370.txt

# STAFFIIDatabase

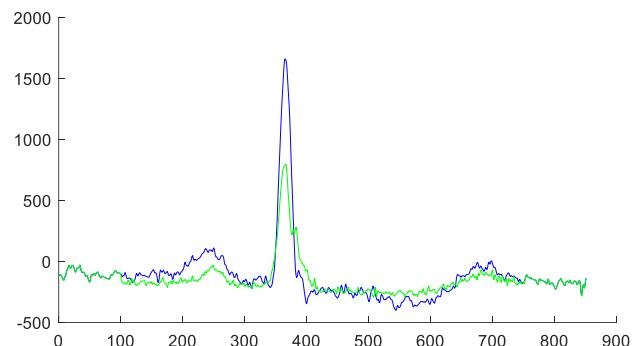
Here we swapped  
in short regions  
from a different  
trace

V6 for V4



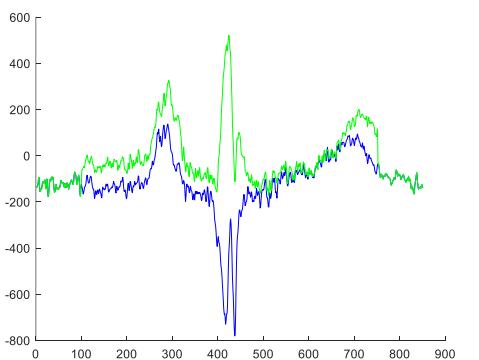
UCR\_Anomaly\_STAFFIIDatabase\_34211\_125720\_126370.txt

V7 for V4



UCR\_Anomaly\_STAFFIIDatabase\_38211\_150720\_151370.txt

8 for 9

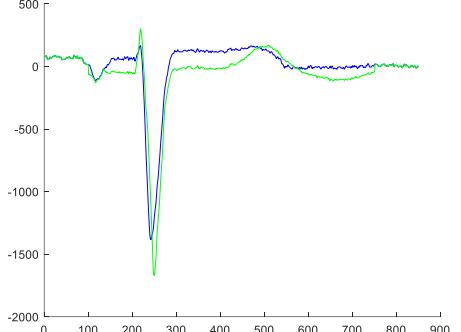


UCR\_Anomaly\_STAFFIIDatabase\_43217\_250720\_251370.txt

## STAFFIIDatabase

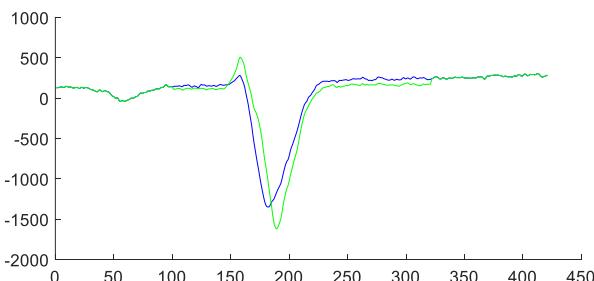
Here we swapped  
in short regions  
from a different  
trace

V2 for V1



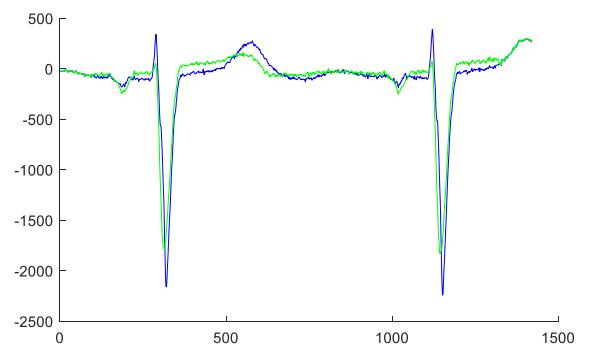
UCR\_Anomaly\_STAFFIIDatabase\_41117\_210720\_211370.txt

V2 for V1



UCR\_Anomaly\_STAFFIIDatabase\_41612\_64632\_64852.txt

V1 for V2



UCR\_Anomaly\_STAFFIIDatabase\_45616\_163632\_164852.txt