

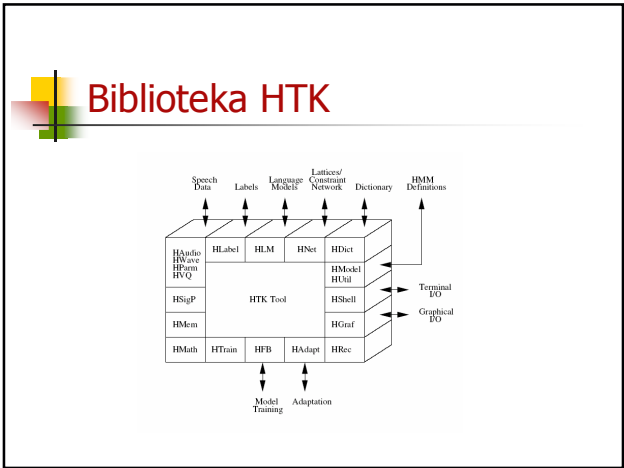


# Biblioteka HTK

---

# Biblioteka HTK

The diagram illustrates the HTK Toolkit architecture. At the center is a 3D block labeled "HTK Tool". Above this block, a row of five boxes represents input/output modules: HAudio, HWave, HForm, HVQ, HLabel, HLM, HNet, HDict, HModel, HUtil, HShell, HGrf, and HRec. Arrows point from these modules to the central HTK Tool block. Above the top row of modules, labels indicate the types of data they handle: Speech Data (for HAudio), Labels (for HLabel), Language MSHell (for HLM), Lattices/Constraint Network (for HNet), Dictionary (for HDict), and HMM Definitions (for HModel). Below the HTK Tool block, a row of five boxes represents processing modules: HMath, HTran, HFB, HAdapt, and HRec. Arrows point from the HTK Tool block to these modules. Below the bottom row of modules, labels indicate the types of processing they perform: Model Training (for HMath), Adaptation (for HAdapt), and Graphical IO (for HGrf). To the right of the HTK Tool block, labels indicate the types of output they produce: Terminal IO (for HModel) and Graphical IO (for HGrf).



# Etapy przetwarzania danych z wykorzystaniem modułów HTK

```

graph TD
    subgraph Data_Prep [Data Prep]
        HLED[HLED]
        HLSTATS[HLSTATS]
        HSLAB[HSLAB]
        HCOPI[HCOPI]
        HLIST[HLIST]
        HQUANT[HQUANT]
    end

    subgraph Training
        Transcriptions[Transcriptions]
        Speech[Speech]
        HCOMPV[HCOMPV, HINT, HREST, HEREST]
        HSMOOTH[HSMOOTH, HHED, HEADAPT]
    end

    subgraph Testing
        HDMAN[HDMAN]
        Dictionary[Dictionary]
        Networks[Networks]
        HVITE[HVITE]
    end

    subgraph Analysis
        HIRILD[HIRILD]
        HPARSE[HPARSE]
        Transcriptions2[Transcriptions]
        HRSULTS[HRSULTS]
    end

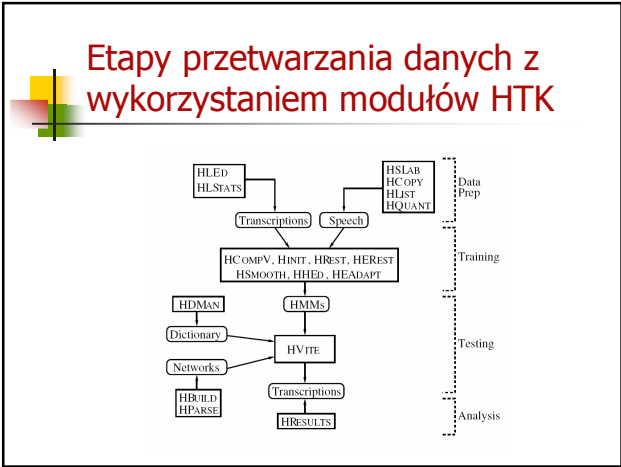
    HLED --> Transcriptions
    HLSTATS --> Transcriptions
    HSLAB --> Speech
    HCOPI --> Speech
    HLIST --> Speech
    HQUANT --> Speech
    Transcriptions --> HCOMPV
    Speech --> HCOMPV
    HCOMPV --> HSMOOTH
    HDMAN --> Dictionary
    Dictionary --> HVITE
    Networks --> HVITE
    HVITE --> Transcriptions2
    Transcriptions2 --> HIRILD
    Transcriptions2 --> HPARSE
    Transcriptions2 --> HRSULTS
  
```

The flowchart illustrates the HTK data processing pipeline, organized into four main stages indicated by dashed brackets on the right:

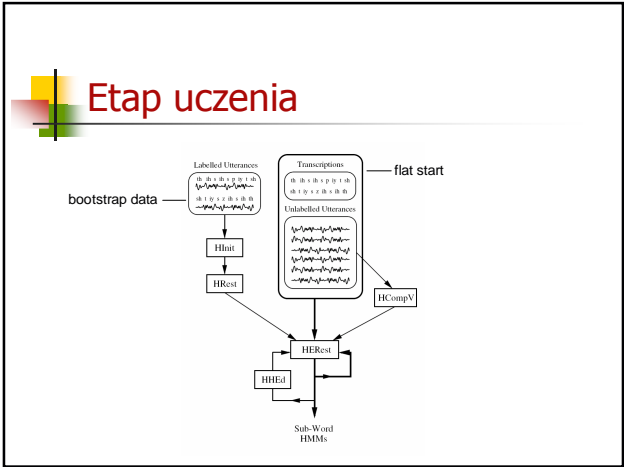
- Data Prep:** Includes modules HLED, HLSTATS, HSLAB, HCOPI, HLIST, and HQUANT.
- Training:** Includes modules Transcriptions, Speech, HCOMPV, HINT, HREST, HEREST, HSMOOTH, HHED, and HEADAPT.
- Testing:** Includes modules HDMAN, Dictionary, Networks, and HVITE.
- Analysis:** Includes modules HIRILD, HPARSE, Transcriptions, and HRSULTS.

The flow of data is as follows:

- HLED and HLSTATS feed into Transcriptions.
- HSLAB, HCOPI, HLIST, and HQUANT feed into Speech.
- Transcriptions and Speech feed into HCOMPV, HINT, HREST, and HEREST.
- HCOMPV, HINT, HREST, and HEREST feed into HSMOOTH, HHED, and HEADAPT.
- HDMAN feeds into Dictionary.
- Dictionary and Networks feed into HVITE.
- HVITE feeds into Transcriptions.
- Transcriptions feeds into HIRILD, HPARSE, and HRSULTS.



The diagram illustrates the HTS algorithm flow. It starts with 'bootstrap data' which is used to create 'Labelled Utterances'. These are processed by 'HInit' and 'HRest' blocks. The 'Transcriptions' (labeled 'flat start') are used to create 'Unlabelled Utterances'. These are processed by 'HcompV' and 'HRest' blocks. The 'HRest' block receives input from both the 'Labelled Utterances' path and the 'Unlabelled Utterances' path. The output of 'HRest' is 'Sub-Word HMMs'.



# Tworzenie systemu rozpoznawania mowy z wykorzystaniem biblioteki **HTK**

## Przygotowanie danych

- krok 1: zdefiniowanie gramatyki

```
graph LR; start([send-start]) --> sms[sms]; start --> chat[chat]; start --> rozmowa[rozmowa]; start --> zadzwon[zadzwoń]; sms --> do[do]; mail[mail] --> do; chat --> z[z]; rozmowa --> z; zadzwon --> 0[0]; zadzwon --> 1[1]; zadzwon --> 9[9]; do --> damska[damska]; do --> michala[michała]; do --> macieja[macieja]; do --> krzysztofa[krzysztofa]; z --> damianem[damianem]; z --> michale[michale]; z --> maciejem[maciejem]; z --> krzysztofem[krzysztofem]; 0 --> end([send-end]); 1 --> end; 9 --> end;
```

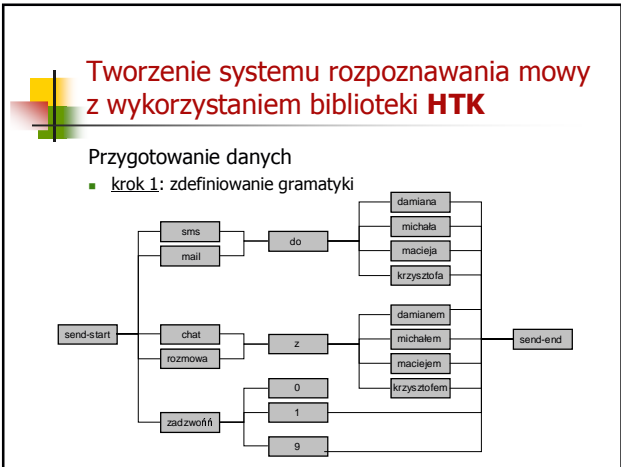
# Tworzenie systemu rozpoznawania mowy z wykorzystaniem biblioteki **HTK**

## Przygotowanie danych

- krok 1: zdefiniowanie gramatyki

```
graph LR; start([send-start]) --> sms[sms]; start --> chat[chat]; start --> rozmowa[rozmowa]; start --> zadzwon[zadzwoń]; sms --> do[do]; mail[mail] --> do; chat --> z[z]; rozmowa --> z; zadzwon --> 0[0]; zadzwon --> 1[1]; zadzwon --> 9[9]; do --> damska[damska]; do --> michala[michała]; do --> macieja[macieja]; do --> krzysztofa[krzysztofa]; z --> damianem[damianem]; z --> michale[michale]; z --> maciejem[maciejem]; z --> krzysztofem[krzysztofem]; 0 --> end([send-end]); 1 --> end; 9 --> end;
```

- # Tworzenie systemu rozpoznawania mowy z wykorzystaniem biblioteki **HTK**
- ## Przygotowanie danych
- krok 1: zdefiniowanie gramatyki
- 
- ```
graph LR; start([send-start]) --> sms[sms]; start --> chat[chat]; start --> rozmowa[rozmowa]; start --> zadzwon[zadzwoń]; sms --> do[do]; mail[mail] --> do; chat --> z[z]; rozmowa --> z; zadzwon --> 0[0]; zadzwon --> 1[1]; zadzwon --> 9[9]; do --> damska[damska]; do --> michala[michała]; do --> macieja[macieja]; do --> krzysztofa[krzysztofa]; z --> damianem[damianem]; z --> michale[michale]; z --> maciejem[maciejem]; z --> krzysztofem[krzysztofem]; 0 --> end([send-end]); 1 --> end; 9 --> end;
```



# Przykład pliku typu *gram*

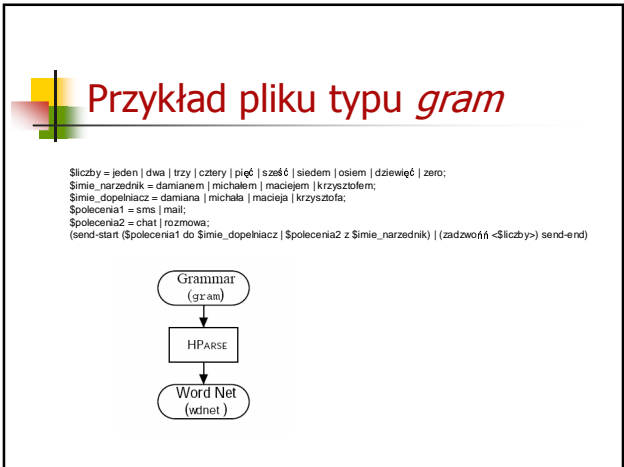
\$liczby = jeden | dwa | trzy | cztery | pięć | sześć | siedem | osiem | dziewięć | zero;  
\$imie\_narzednik = damianem | michelem | maciejem | krzysztofem;  
\$imie\_dopelniacz = damiana | micheła | macieja | krzysztofa;  
\$spolecenia1 = sms | mail;  
\$spolecenia2 = chat | rozmowa;  
(send-start (\$spolecenia1 do \$imie\_dopelniacz | \$spolecenia2 z \$imie\_narzednik) | (zadzwoń! <\$liczby>) send-end)

```
graph TD; A([Grammar  
(gram)]) --> B[HPARSE]; B --> C([Word Net  
(wtnet)])
```

# Przykład pliku typu *gram*

\$liczby = jeden | dwa | trzy | cztery | pięć | sześć | siedem | osiem | dziewięć | zero;  
\$imie\_narzednik = damianem | michelem | maciejem | krzysztofem;  
\$imie\_dopelniacz = damiana | micheła | macieja | krzysztofa;  
\$spolecenia1 = sms | mail;  
\$spolecenia2 = chat | rozmowa;  
(send-start (\$spolecenia1 do \$imie\_dopelniacz | \$spolecenia2 z \$imie\_narzednik) | (zadzwoń! <\$liczby>) send-end)

```
graph TD; A([Grammar  
(gram)]) --> B[HPARSE]; B --> C([Word Net  
(wtnet)])
```





## przykład pliku typu *wdnet*

```

VERSION=1.0
N=33 L=61
l=0 W=send-end
l=1 W=NULL
l=2 W=zero
l=3 W=NULL
l=4 W=dziewi352/346
l=5 W=osiem
l=6 W=siedem
l=7 W=pi352/346
l=8 W=pi352/346
l=9 W=cztery
l=10 W=trzy
l=11 W=dwa
l=12 W=jeden
l=13 W=zadzwon
l=14 W=krzysztofem
l=15 W=maciejem
l=16 W=michal263em
l=17 W=damianem
l=18 W=
l=19 W=rozmowa
l=20 W=NULL
l=21 W=chat

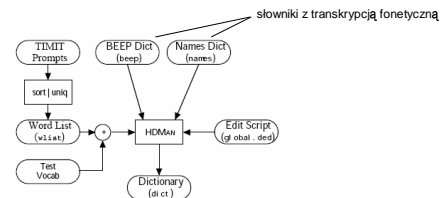
l=22 W=krzysztofa
l=23 W=macieja
l=24 W=michal263a
l=25 W=damiana
l=26 W=do
l=27 W=mail
l=28 W=NULL
l=29 W=sms
l=30 W=send-start
l=31 W=NULL
l=32 W=NULL
l=33 W=NULL

J=13 S=4 E=3
J=14 S=5 E=3
J=15 S=6 E=3
J=16 S=7 E=3
J=17 S=8 E=3
J=18 S=9 E=3
J=19 S=10 E=3
J=20 S=11 E=3
J=21 S=12 E=3
J=22 S=3 E=4
J=23 S=13 E=4
J=24 S=3 E=5
J=25 S=13 E=5
J=26 S=3 E=6
J=27 S=13 E=6
J=28 S=3 E=7
J=29 S=13 E=7
J=30 S=3 E=8
J=31 S=13 E=8
J=32 S=3 E=9
J=33 S=13 E=9
J=34 S=3 E=10
J=35 S=13 E=10
J=36 S=3 E=11

J=37 S=13 E=11
J=38 S=3 E=12
J=39 S=13 E=12
J=40 S=32 E=13
J=41 S=18 E=14
J=42 S=18 E=15
J=43 S=18 E=16
J=44 S=18 E=17
J=45 S=20 E=18
J=46 S=30 E=19
J=47 S=19 E=20
J=48 S=21 E=20
J=49 S=30 E=21
J=50 S=26 E=22
J=51 S=26 E=23
J=52 S=26 E=24
J=53 S=26 E=25
J=54 S=26 E=26
J=55 S=30 E=27
J=56 S=27 E=28
J=57 S=29 E=28
J=58 S=30 E=29
J=59 S=32 E=30
J=60 S=1 E=31
    
```

## Moduł *HDMan*

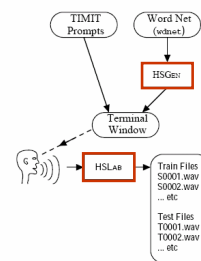
- krok 2: słownik z transkrypcją fonetyczną
  - utworzenie posortowanej listy słów – *wlist*
  - zastosowanie modułu *HDMan* do wygenerowania słownika



## Przykładowy plik *dict*

|          |                       |               |                           |
|----------|-----------------------|---------------|---------------------------|
| jeden    | j e d e n s p         | macieja       | m a c i e j a s p         |
| dwa      | d w a s p             | maciejem      | m a c i e j e m s p       |
| trzy     | t r z y s p           | krzysztofa    | k r z y s z t o f a s p   |
| cztery   | c z e r y s p         | krzysztofem   | k r z y s z t o f e m s p |
| pięć     | p j e n c s p         | chat          | c h a t s p               |
| sześć    | s z e s c s p         | z             | z s p                     |
| siedem   | s i e d e m s p       | do            | d o s p                   |
| osiem    | o s i e m s p         | rozmowa       | r o z m o w a s p         |
| dziwięć  | Z @ e v j e n Z @ s p | sms           | e s m e z s p             |
| zero     | z e r o s p           | mail          | m e j l s p               |
| zadzwon  | z a z v o n s p       | send-end []   | s i l                     |
| wywołaj  | v y v o a j s p       | send-start [] | s i l                     |
| damiana  | d a m j a n a s p     |               |                           |
| damianem | d a m j a n e m s p   |               |                           |
| michala  | m i x a w a s p       |               |                           |
| michalem | m i x a w e m s p     |               |                           |

## Rejestracja i etykietowanie danych akustycznych



## *HSGen*

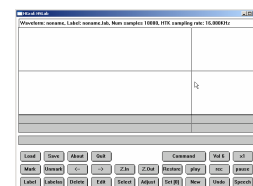
- Przygotowanie pliku z *promptami*:  
HSGen -l 18 wdnet dict > testprompts

### Przykładowy plik *testprompts*:

- chat z maciejem
- zadzwon pięć cztery trzy pięć trzy sześć jeden trzy sześć jeden trzy sześć sześć cztery zero siedem trzy
- chat z maciejem
- chat z micalem
- zadzwon jeden sześć trzy dwa
- zadzwon pięć zero zero pięć trzy osiem dziewięć trzy dziewięć sześć dwa jeden jeden dwa
- chat z krzysztofem
- zadzwon sześć siedem cztery dziewięć dwa dwa jeden cztery pięć zero sześć zero
- zadzwon osiem dwa siedem cztery dziewięć dziewięć cztery jeden
- rozmowa z maciejem
- chat z maciejem
- sms do damiana
- mail do krzysztofa
- mail do michala
- zadzwon osiem
- zadzwon pięć trzy dziewięć dwa zero dwa dwa zero
- zadzwon siedem siedem pięć dwa dwa osiem sześć jeden trzy zero
- zadzwon osiem siedem siedem osiem dwa dziewięć dwa jeden zero sześć

## Moduł *HSLab*

- krok 3: rejestracja i oznaczenie nagrań akustycznych



output:  
noname\_0  
noname\_1

prompt batch

1.wav  
2.wav  
...

### przykład pliku 'prompt batch':

```

import os
import sys
import shutil

output = ['noname_0', 'noname_1']
path = "C:\\\\noname_0"
text = file(path+"testprompts", "r").readlines()
id = 1
i = 0
for s in text:
    print s
    c = shutil.getsize(s)
    if c == 'q':
        break
    os.rename(path+output[i], path+str(id)+".wav")
    i = (i+1) % 2
    id += 1
    
```



## Transkrypcja

### krok 4: generowanie plików z transkrypcją

1. chat z maciejem
2. zadzwonił pięć cztery trzy pięć trzy sześć jeden trzy sześć trzy sześć sześć cztery zero siedem trzy
3. ...

transkrypcja ortograficzna (prompts2mlf)

```
#!MLF#
"1.lab"
chat
z
maciejem
.
"2.lab"
zadzwonił
pięć
cztery
trzy
pięć
trzy
...
```

```
sześć
jeden
trzy
sześć
trzy
sześć
zadzwonił
pięć
cztery
siedem
trzy
.
"3.lab"
...
```

## transkrypcja cd.

words.mlf

```
#!MLF#
sześć
jeden
chat
z
maciejem
.
"2.lab"
zadzwonił
pięć
cztery
trzy
pięć
trzy
.
"3.lab"
...
```

transkrypcja fonetyczna (HLEd)

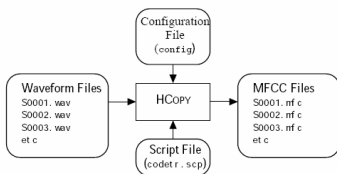
words.mlf

phones0.mlf

```
#!MLF#
p
s6
s6
t
"1.lab"
e
e
e
s6
c6
\361
\346
\346
\346
sil
"3.lab"
...
a
\346
j
d
t
m
e
e
e
a
r
n
c6
\346
y
t
t
e
j
s6
y
r
e
y
s6
y
p
e
z
sil
\234
e
\346
r
.
"2.lab"
\361
t
o
\346
s6
\234
z
t
e
a
s6
d
v
o
\361
```

## Kodowanie – moduł HCOPY

### krok 5: ekstrakcja cech akustycznych z plików wav



### Plik konfiguracyjny (config) modułu HCOPY

```
TARGETKIND = MFCC_0
TARGETRATE = 100000
SAVECOMPRESSED = TRUE
SAVEWITHCRC = TRUE
WINDOWSIZE = 250000.0
USEHAMMING = TRUE
PREEMCOEF = 0.97
NUMCHANS = 26
CEPLIFTER = 22
NUMCEPS = 12
ENORMALISE = FALSE
```

Mel Freq Cepstral Coefficients:

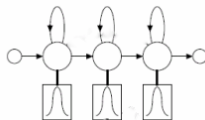
$$c_k = \sqrt{\frac{1}{N} \sum_{n=0}^{N-1} \left( \frac{1}{N} \sum_{m=0}^{N-1} x(m) e^{-j \frac{2\pi}{N} km} \right)^2}$$

MFCC filterbank

## Tworzenie modeli HMM dla monofonów

### krok 6: tworzenie prototypowego modelu – proto (definiowanie topologii modelu, np. 3-stanowe modele HMM)

```
<?xml?>
<HMM?>
  <name?>proto</name?>
  <states?>3</states?>
  <transitions?>3</transitions?>
  <initial?>1</initial?>
  <final?>3</final?>
  <parameters?>
    <state?>1</state?>
    <state?>2</state?>
    <state?>3</state?>
  </parameters?>
</HMM?>
```



## Moduł HCompV

Użycie:

```
HCompV -C config -f 0.01 -m -S train.scf - M hmm0 proto
```

```
<?xml?>
<HMM?>
  <name?>proto</name?>
  <states?>3</states?>
  <transitions?>3</transitions?>
  <initial?>1</initial?>
  <final?>3</final?>
  <parameters?>
    <state?>1</state?>
    <state?>2</state?>
    <state?>3</state?>
  </parameters?>
</HMM?>
```



# Re-estymacja modeli *hmm0*

Użycie 3x:

```
HERest -C config -I phones0.mlf -t 250.0 150.0 1000.0 \
-S train.scp -H hmm0/macros -H hmm0/hmmdefs -M hmm1/monophones0
```

pruning parameters:

macros:  
-o dAFCC\_0\_D\_A-  
-v refSize 30  
-> varFloor1  
-v instance 30 7.320056e-001 2.330221e-001 4.921546e-001 2.679774e-001 ...

```
graph TD
    A([Prototype HMM Definition  
(pr o s)]) --> B[HCompV]
    B --> C([hmm0  
mcr os  
hmdl s])
    D([Training Files  
listed in  
(t r a i n. s c p)]) --> E[HERest]
    F([HMM list  
( monophones0 )]) --> E
    C --> E
    E --> G([hmm1  
mcr os  
hmdl s])
    H([Phone Level  
Transcription  
(phones0.mlf)]) --> E
```

# Wynik re-estymacji hmm\*

## Poprawienie modeli ciszy (silence models) poprzez dodanie dodatkowego stanu

- krok 7: dodanie dodatkowego przejścia między stanem 2 i 4 oraz 4 i 2 w modelu ciszy oraz dodanie stanu krótkiej pauzy *sp*

The diagram illustrates a Hidden Markov Model (HMM) for silence. It consists of three states: 'sil' (state 2), a central state, and '4'. Transitions are shown between these states. A red circle highlights a sub-model with states 'sp' and 'shared state', which is connected to the main model.

## Krok 7 - wykonanie

- skopiować środkowy stan modelu ciszy *sil*, aby utworzyć nowy model *sp* i zapisać wynik w pliku `hmmdefs` czwartej iteracji (katalog `hmm4`),
- uruchomić moduł edytora HMM - HHed, aby dodać dodatkowe przejścia między stanami oraz aby podłączyć model *sp* do centralnego stanu modelu *sil*

```
HHed -H hmm4/macros -H hmm4/hmmdefs -M hmm5 sil.hed monophones1
```

```
sil.hed:
AT 2 4 0.2 {sil.transp}
AT 4 2 0.2 {sil.transp}
AT 1 3 0.3 {sp.transp}
TI silst {sil.state[3], sp.state[2]}
```

- uruchomić dwukrotnie moduł re-estymacji - HERest

## Podsumowanie kroku 7

```

graph LR
    A[Edit sil -> sp] --> B(hmm4 mcr os hmdef s)
    C(HMM list monophones1) --> D(HHEd)
    E(Edit Script sil.hed) --> D
    B --> D
    D --> F(hmm5 mcr os hmdef s)
    F --> G(HERest x2)
    G --> H(hmm7 mcr os hmdef s)
  
```

# Uwzględnienie różnej wymowy określonych słów – moduł HVite

- krok 8:**
  - uruchomić moduł HVite:

```
HVite -l '*' -o SWT -b silence -C config -a -H hmm7/macros \
-H hmm7/hmmdefs -i aligned.mlf -m -t 250.0 -y lab \
-I words.mlf -S train.scp dict monophones1
```

```

graph TD
    WLT[Word Level Transcriptions  
(words.mlf)] --> HVite
    D[Dictionary  
(dict)] --> HVite
    HML[HMM list  
(monophones1)] --> HVite
    TD[Timing Data  
Labeled in  
(train.scp)] --> HVite
    HVite --> FB[formos  
bundle.s]
    HVite --> PLT[Phone Level  
Transcriptions  
(all good.mlf)]
  
```

- Ponownie uruchomić dwukrotnie moduł re-estymacji – HERest -> hmm9



## Multifony

- **krok 9: tworzenie (di)trifonów z monofonów (word internal)**

```
HLEd -n triphones1 -l '*' -i wintri.mlf mktri.led aligned.mlf
```

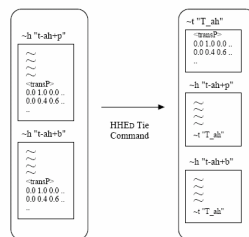
| mktri.led: | aligned.mlf: | wintri.mlf: | + | triphones           |
|------------|--------------|-------------|---|---------------------|
| WB sp      | #MLF!#       | #MLF!#      |   | (wszystkie możliwe) |
| WB sil     | "t1.lab"     | "t1.lab"    |   |                     |
| TC         | sil          | sil         |   |                     |
|            | c6           | c6+a        |   |                     |
|            | a            | c6-a+d      |   |                     |
|            | d            | a-d         |   |                     |
|            | sp           | z           |   |                     |
|            | z            | sp          |   |                     |
|            | m            | m+a         |   |                     |
|            | a            | m-a\346     |   |                     |
|            | \346         | a\346+e     |   |                     |
|            | e            | \346-e+j    |   |                     |
|            | j            | e+e         |   |                     |
|            | e            | e+m         |   |                     |
|            | m            | e-m         |   |                     |
|            | sp           | sp          |   |                     |
|            | sil          | sil         |   |                     |

## Tworzenie modeli multifonów

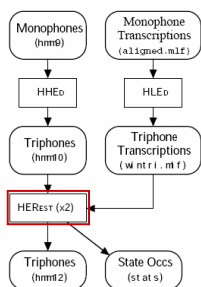
```
HHEd -B -H hmm9/macros -H hmm9/hmmdefs -M hmm10 mktri.hed monophones1
```

mktri.hed:

```
CL triphones1
TI T_Z @ {("Z@+","Z@+","Z@+").transP}
TI T_J {("J+","J+","J+").transP}
TI T_C @ {("C@+","C@+","C@+").transP}
TI T_X {("X+","X+","X+").transP}
TI T_Z {("Z+","Z+","Z+").transP}
TI T_d {("d+","d+","d+").transP}
TI T_a {("a+","a+","a+").transP}
TI T_e {("e+","e+","e+").transP}
TI T_d {("d+","d+","d+").transP}
TI T_g {("g+","g+","g+").transP}
TI T_6 {("6+","6+","6+").transP}
TI T_j {("j+","j+","j+").transP}
TI T_k {("k+","k+","k+").transP}
TI T_j {("j+","j+","j+").transP}
TI T_m {("m+","m+","m+").transP}
TI T_s6 {("s6+","s6+","s6+").transP}
TI T_o {("o+","o+","o+").transP}
TI T_n {("n+","n+","n+").transP}
```



## Re-estymacja w kroku 9



## Łączenie stanów modeli trifonów

- **krok 10:**

```
HHEd -B -H hmm12/macros -H hmm12/hmmdefs -M hmm13 tree.hed triphones1
```

