

Automatic Music Score Recognition

Oral Defense
HW Lee, 07.11.2016

Optical Music Recognition

- Score in Western Music
- Image to Symbolic Representation (MusicXML)



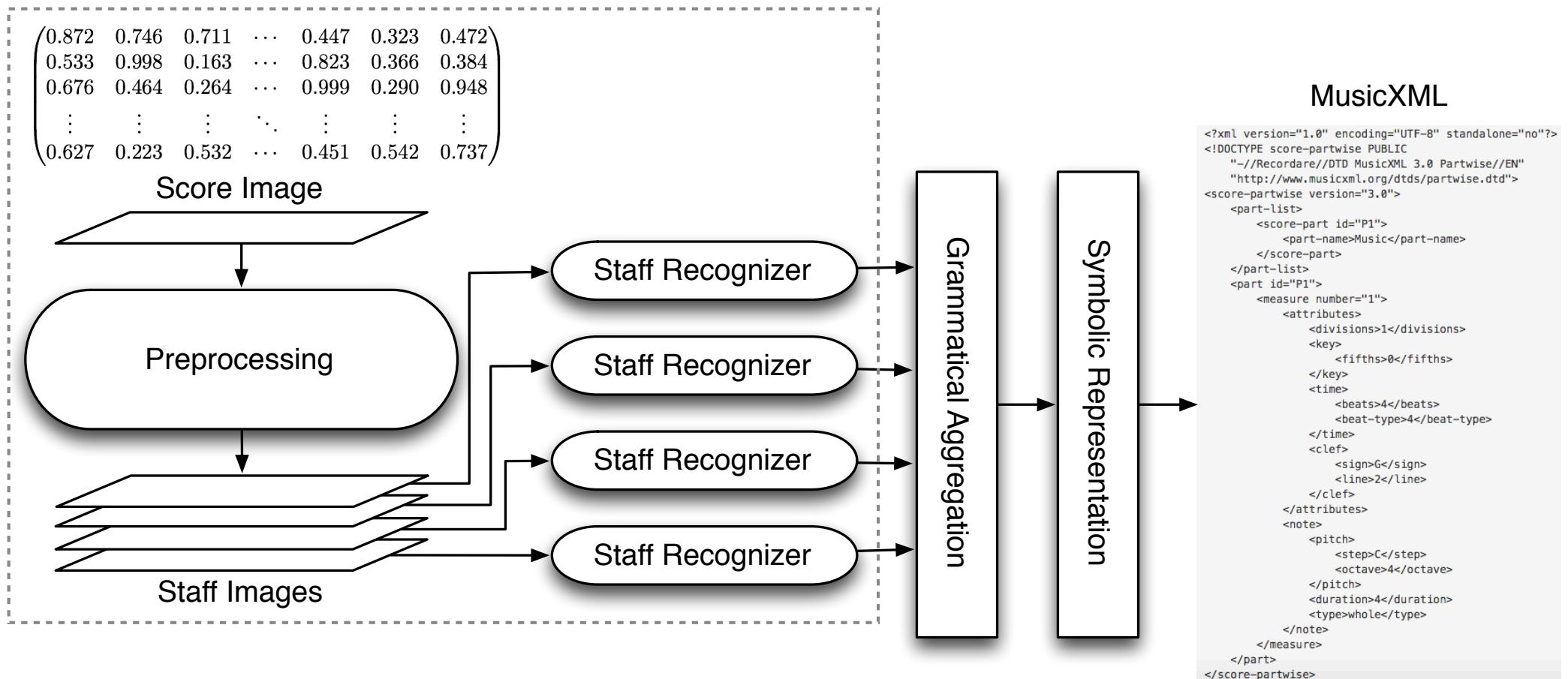
For human

0.872	0.746	0.711	...	0.447	0.323	0.472
0.533	0.998	0.163	...	0.823	0.366	0.384
0.676	0.464	0.264	...	0.999	0.290	0.948
:	:	:	..	:	:	:
0.627	0.223	0.532	...	0.451	0.542	0.737

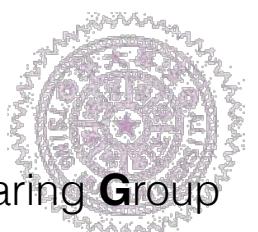
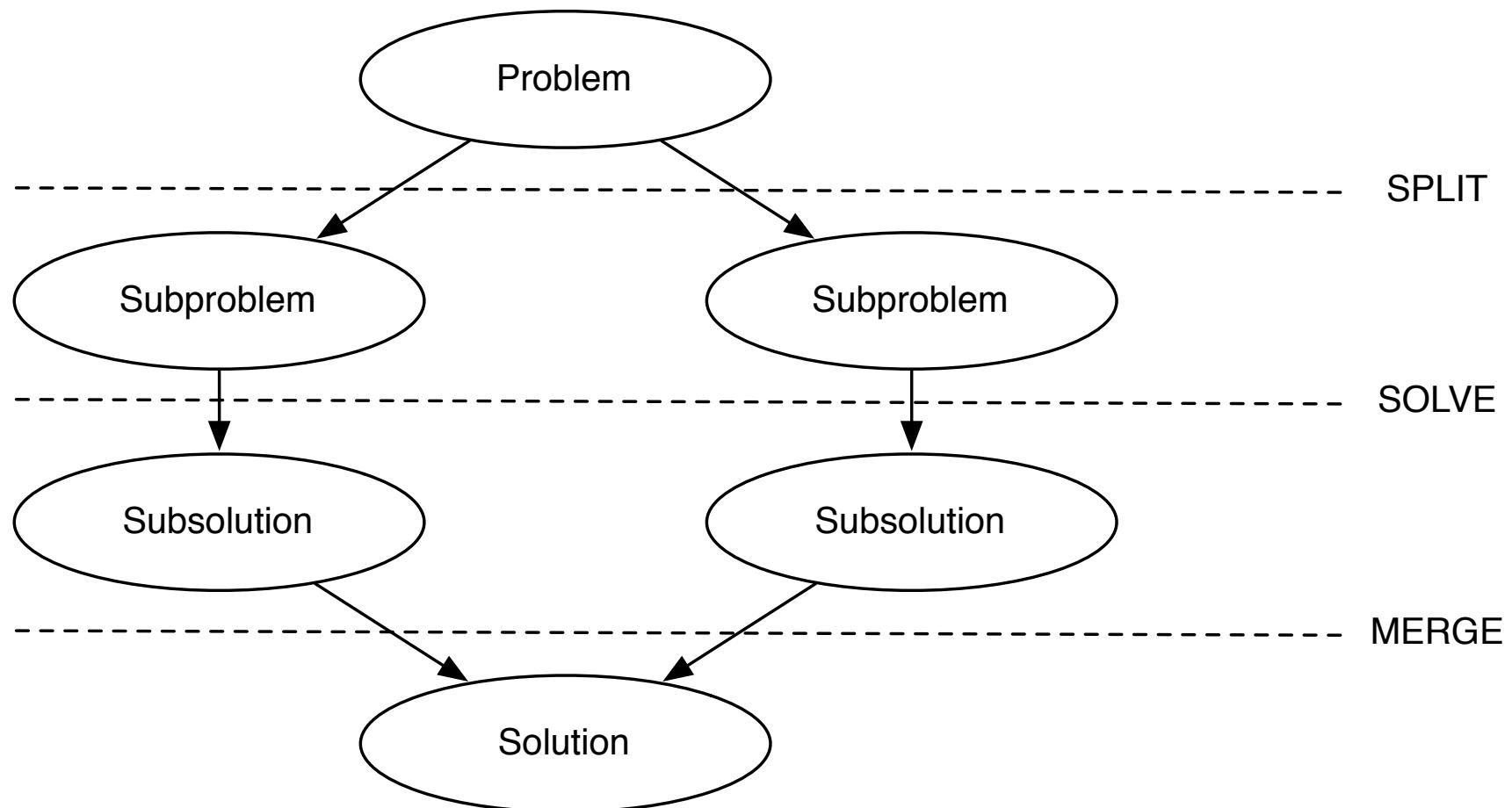
For computer



A Typical OMR System



Divide & Conquer



Score Sample

好久不見
陳奕迅

詞 施立
曲 陳小霞
編曲 大騷

$\text{♩} = 68$

Solo
Soprano
Alto
Tenor
Bass

我來到一
oo oo oo oo
oo oo oo oo
oo oo oo oo
oo oo oo oo
oo dm

好想你

作詞：黃明志
作曲：黃明志
編曲：陳爌安

$\text{♩} = 96$

俐安
仲文
HW
PK
ALN

想要 傳送一封簡訊給你 我好想好想你 想要
ba ba ba du wa
ba ba ba du wa
ba ba ba du wa
dm dm dm dm dm

6

Solo
S.
A.
T.
B.

你的城市一 走過你來時一的路 想像著一沒
dm dm dm du

第八格人聲樂團 ©

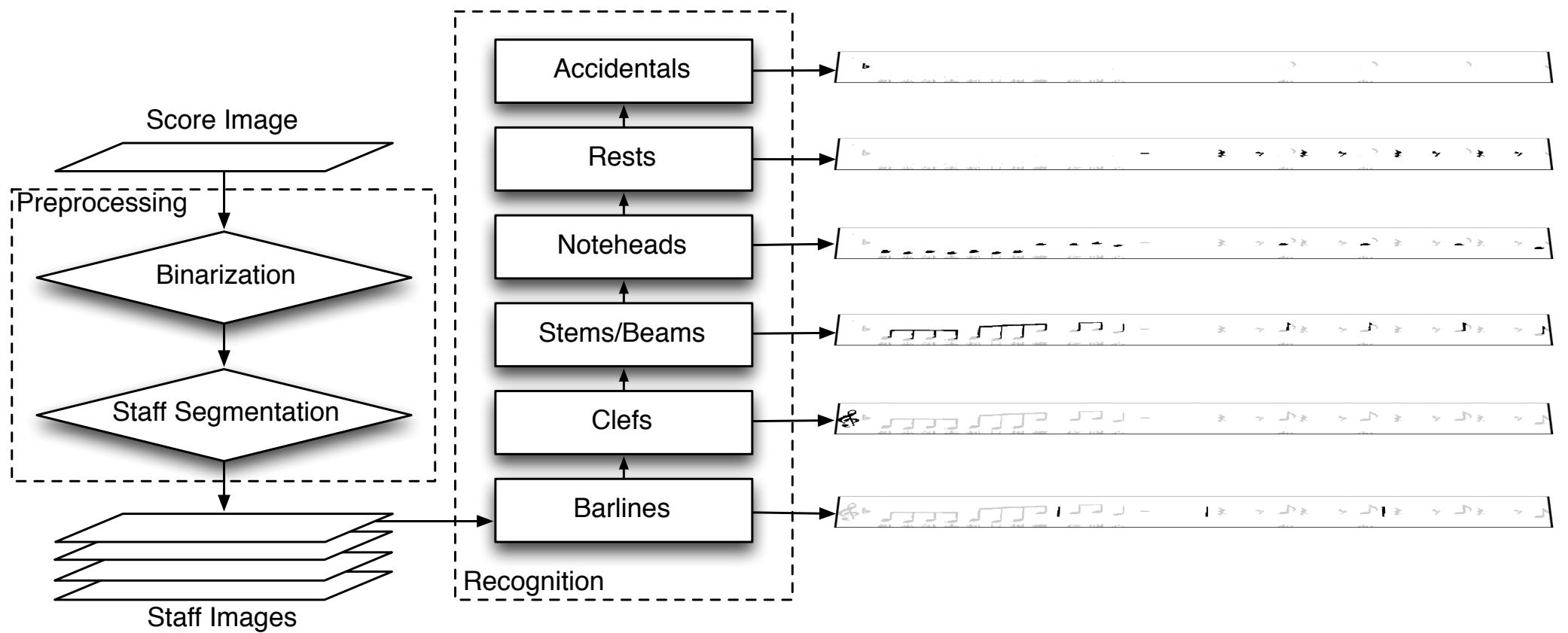
S.
A.
T.
T.
B.

立刻打通電話給你 我好想好想你 ba ba
ba ba du du du 每天起床的第一件事情 就是
ba ba du du du ba ba
dm dm dm dm dm

第八格人聲樂團 ©

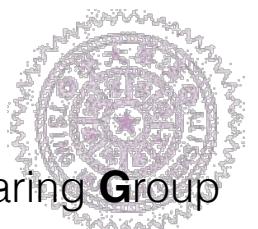
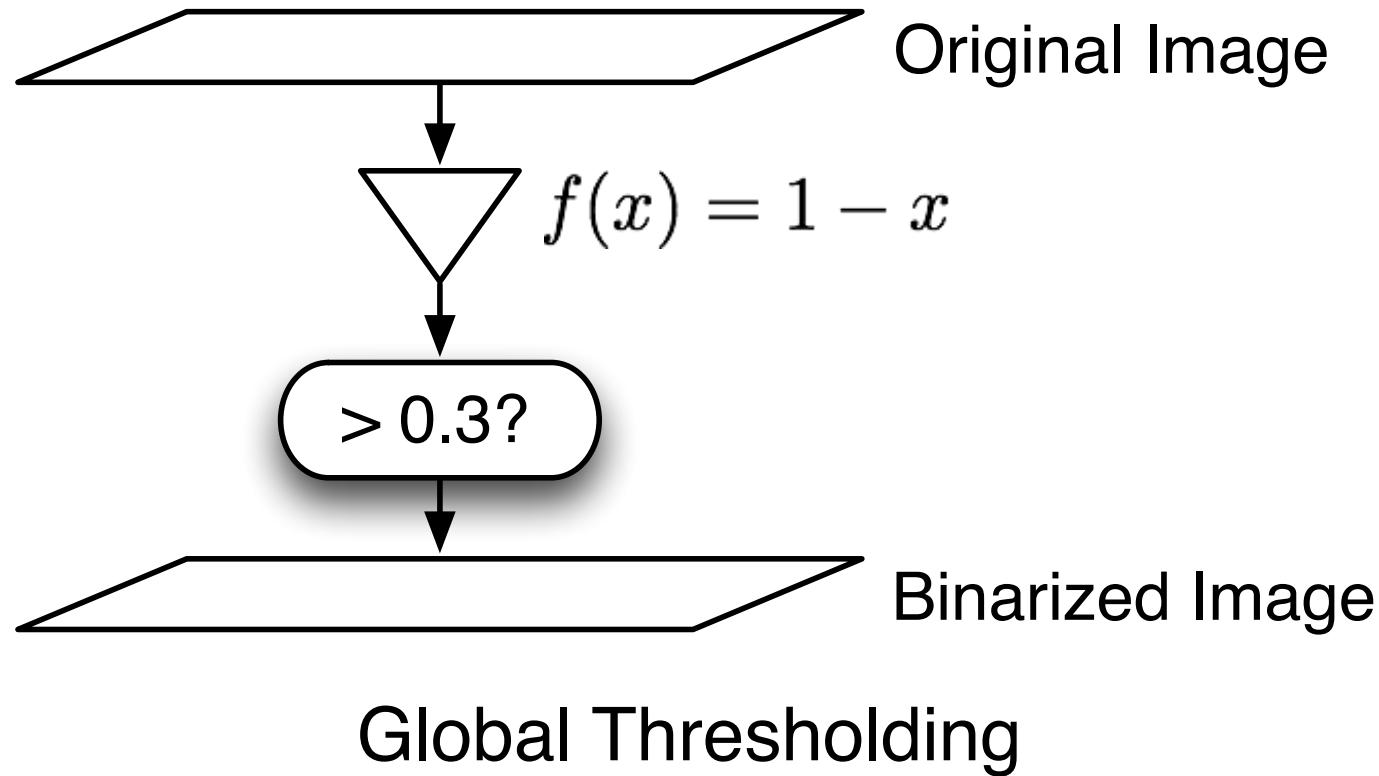


System Overview



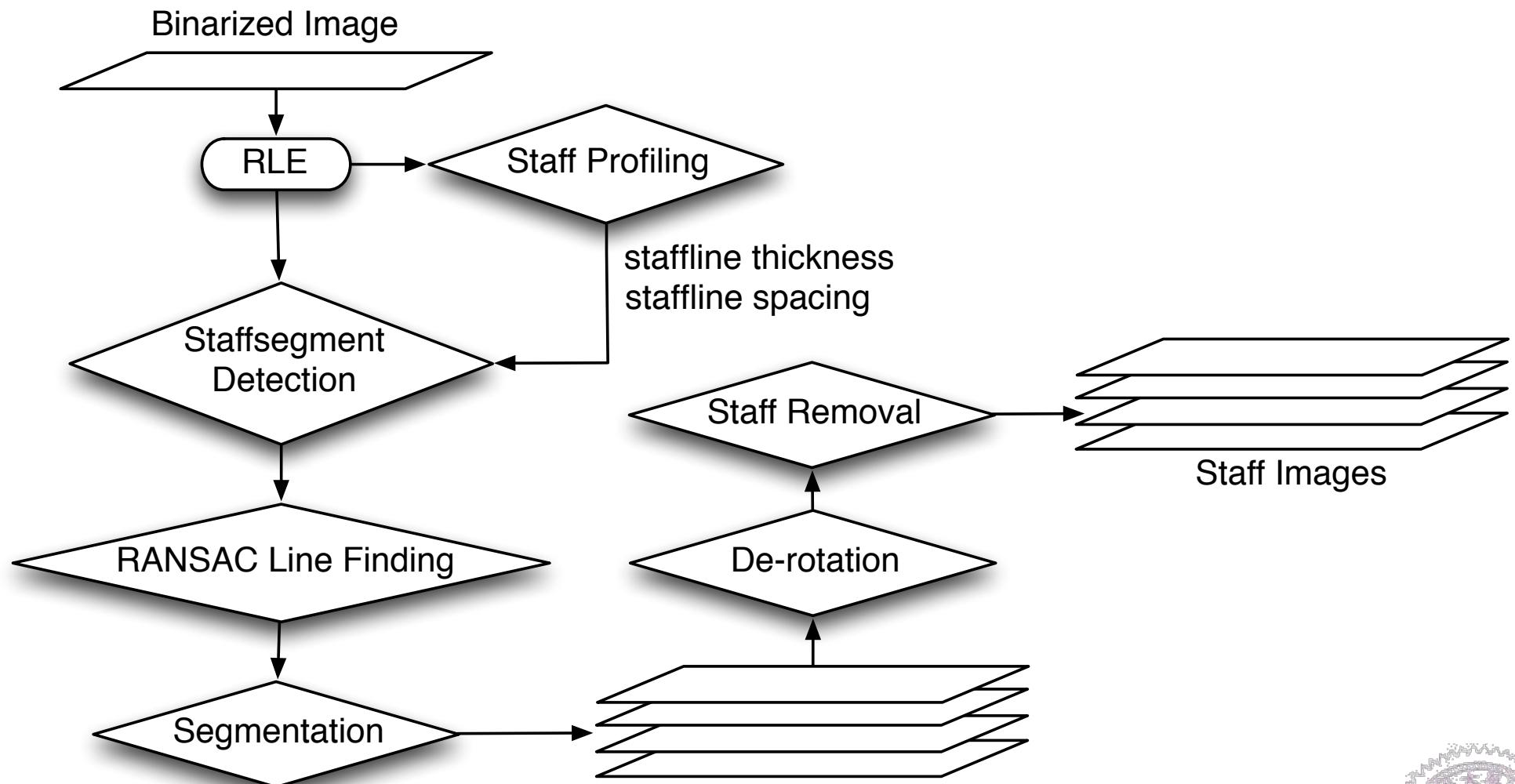
Preprocessing Binarization

Binarization

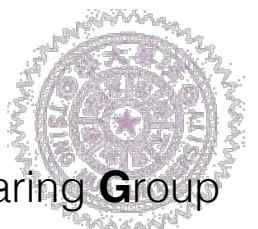
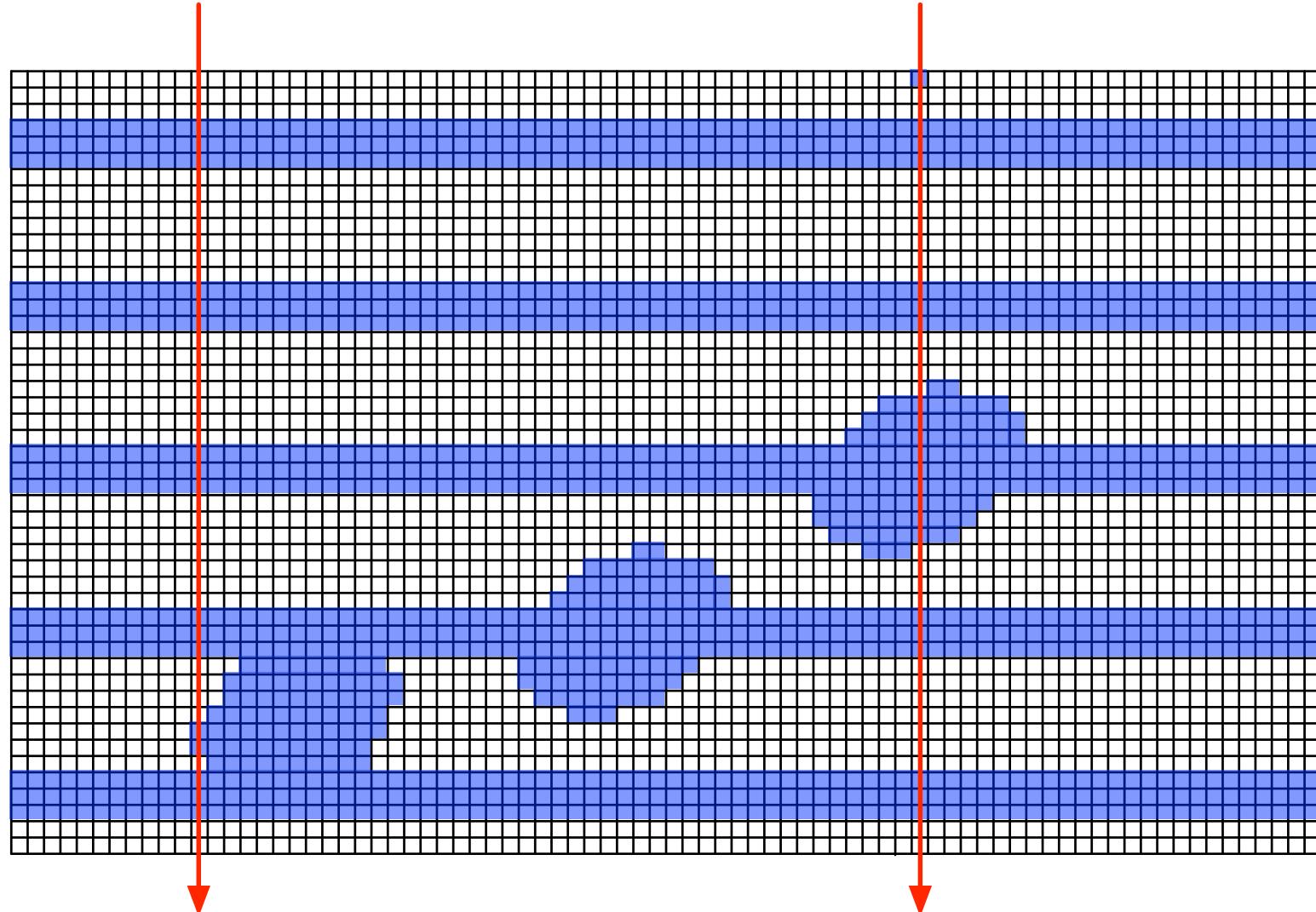


Preprocessing Staff Segmentation

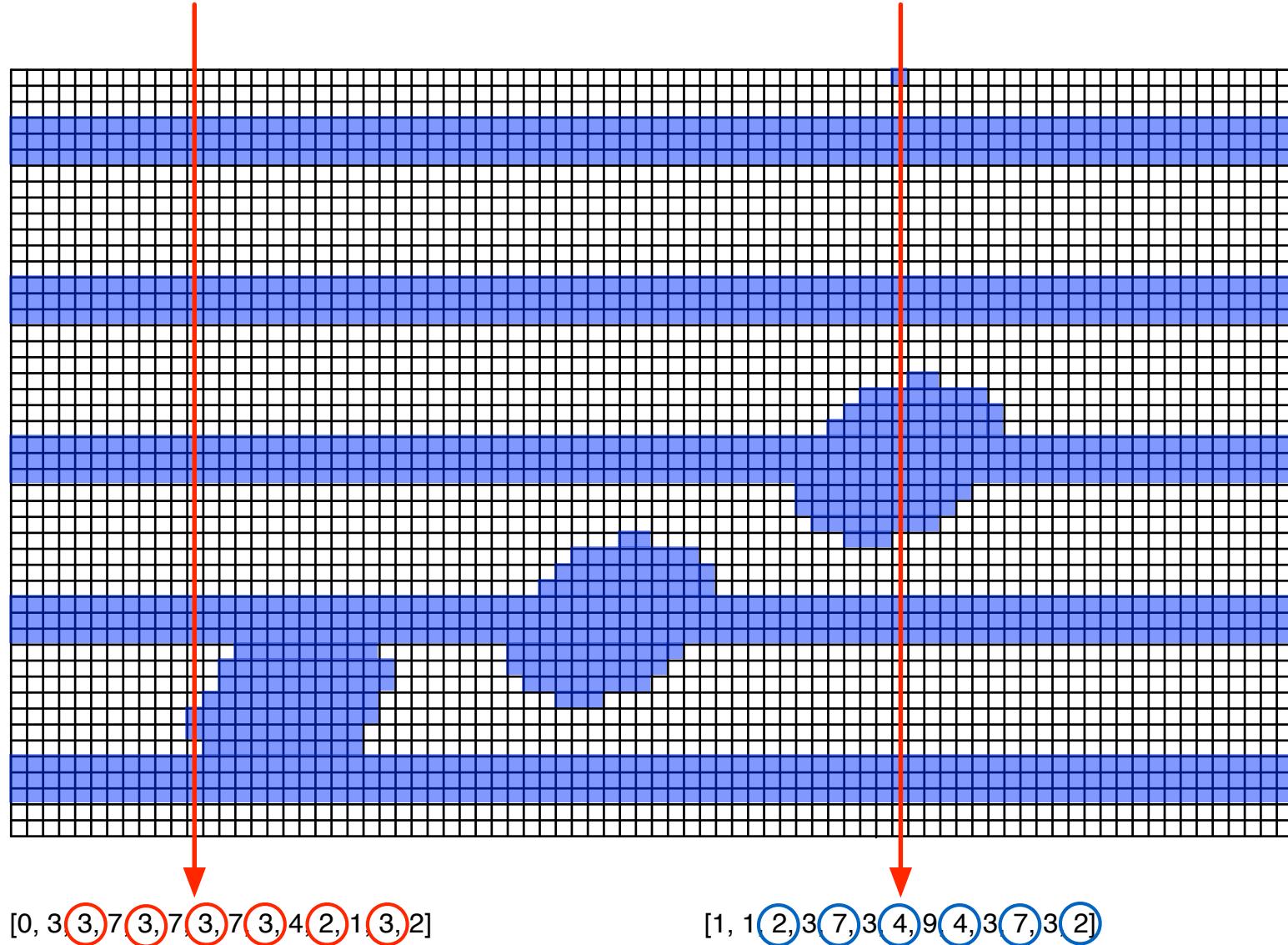
Staff Segmentation



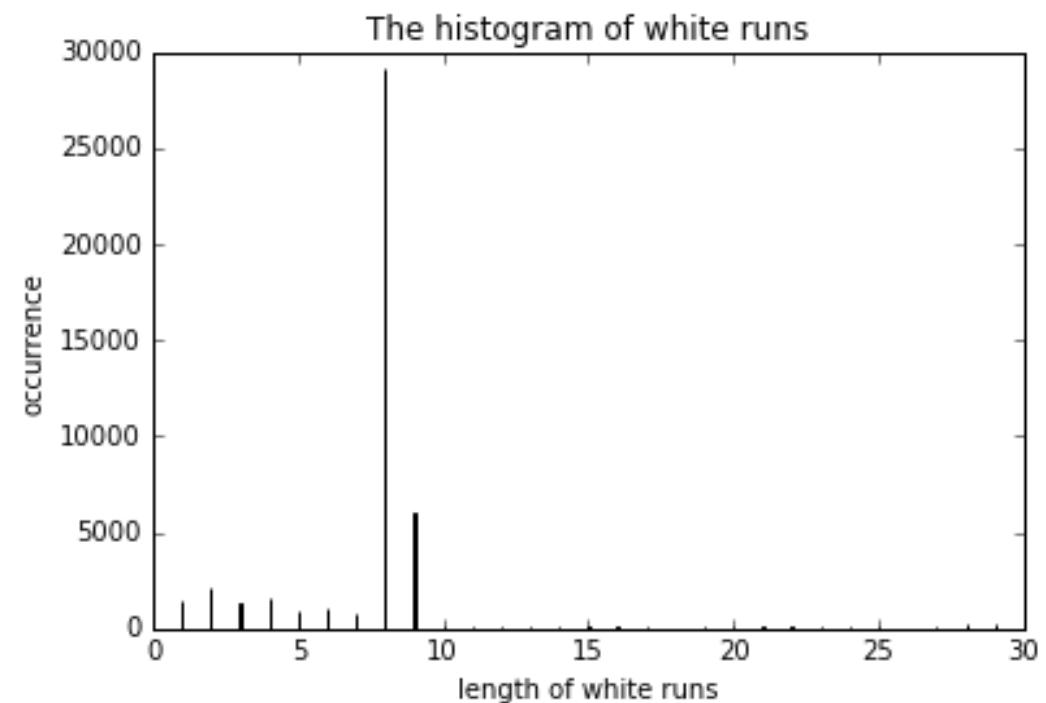
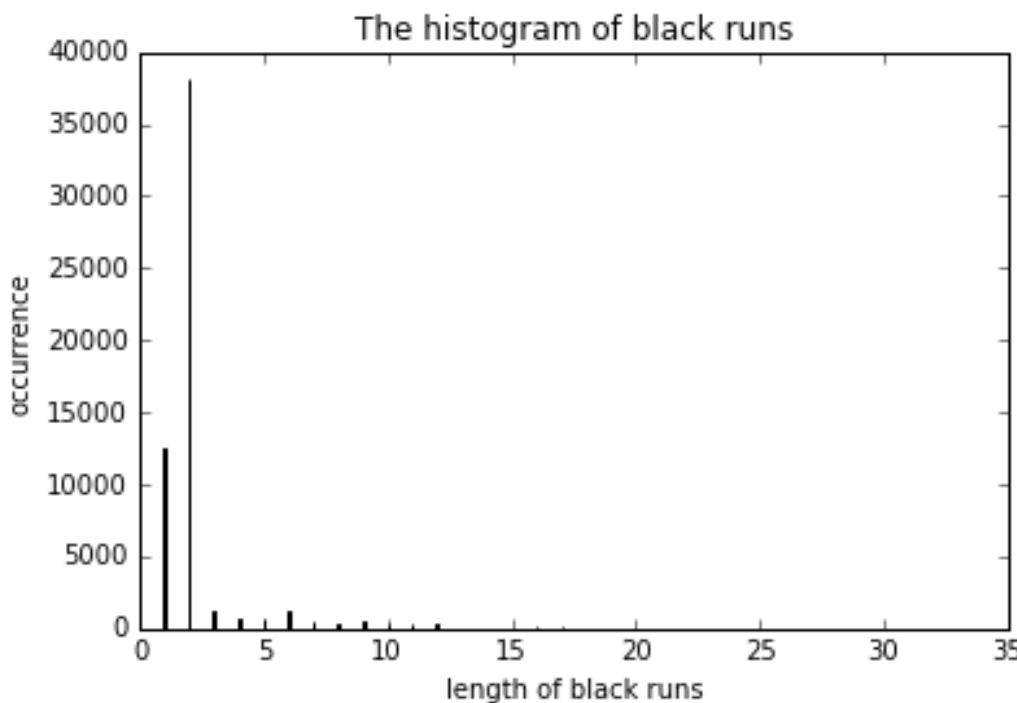
RLE (Run-length Encoding)



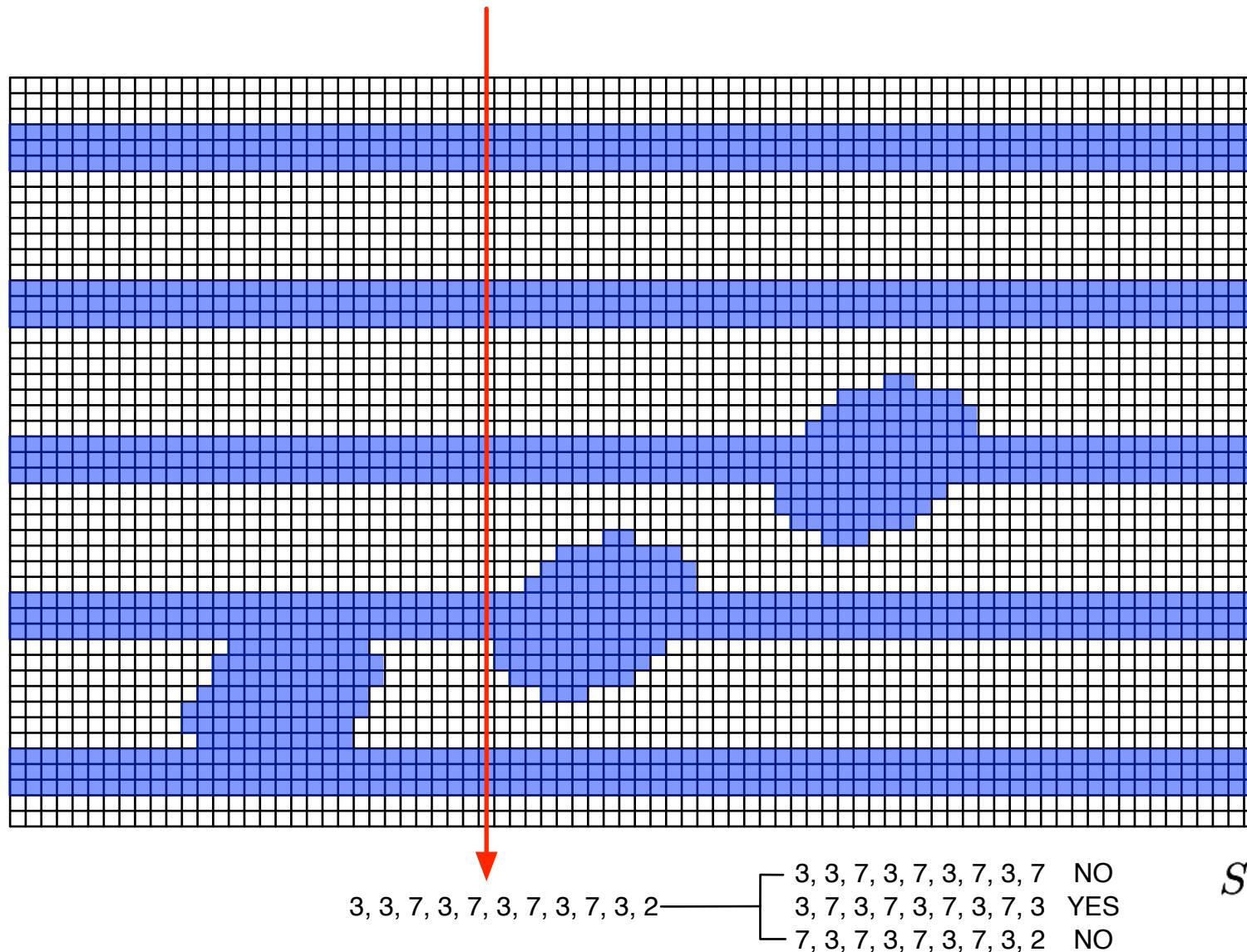
Staff Profile



Staff Profile

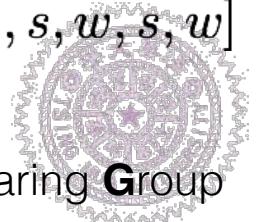


Staffsegment Detection



$$e(v') = \frac{\|v' - S\|_1}{\|S\|_1},$$

$$S = [w, s, w, s, w, s, w, s, w]$$



Staffsegment

40

S. 說來說去都只想讓 你開心 好想 你 好想你 好想 你 好想 你 是

A. 說來說去都只想讓 你開心 du du

T. 說來說去都只想讓 你開心 du

T. oo oo du

B. dm dm

The musical score consists of five staves, each representing a different vocal part: Soprano (S.), Alto (A.), Tenor (T.), Bass (T.), and Bassoon (B.). The score is set in common time (indicated by '40'). The lyrics are written in Chinese characters below each staff. The music features a mix of eighth and sixteenth notes, with some sustained notes indicated by red horizontal lines. The bassoon part (B.) has a unique bass clef and note heads.



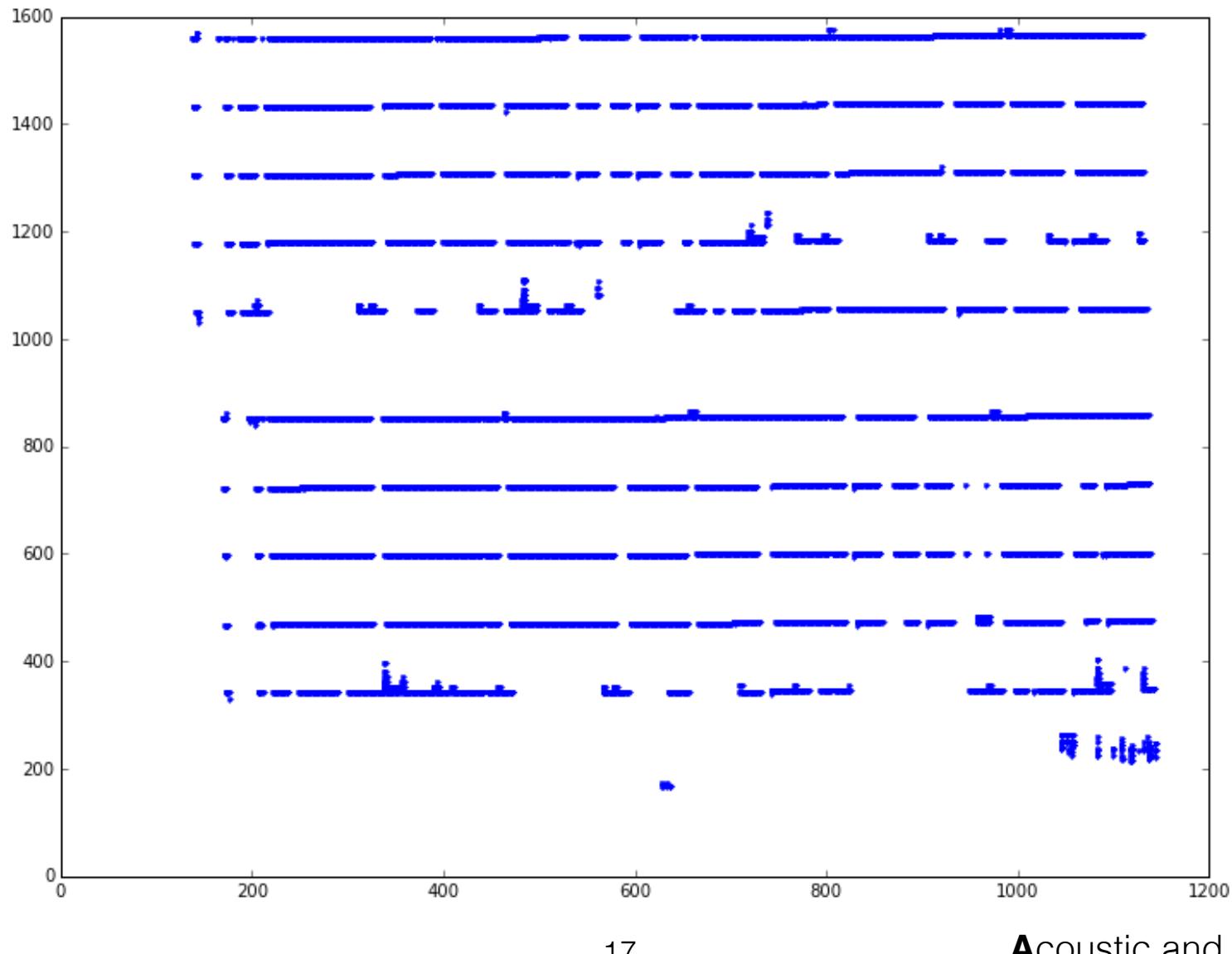
Staffsegment

A musical score for five voices (Soprano, Alto, Tenor, Bass, and a fifth tenor) on five staves. The key signature is A minor (no sharps or flats). The time signature is common time (indicated by '4'). The vocal parts are:

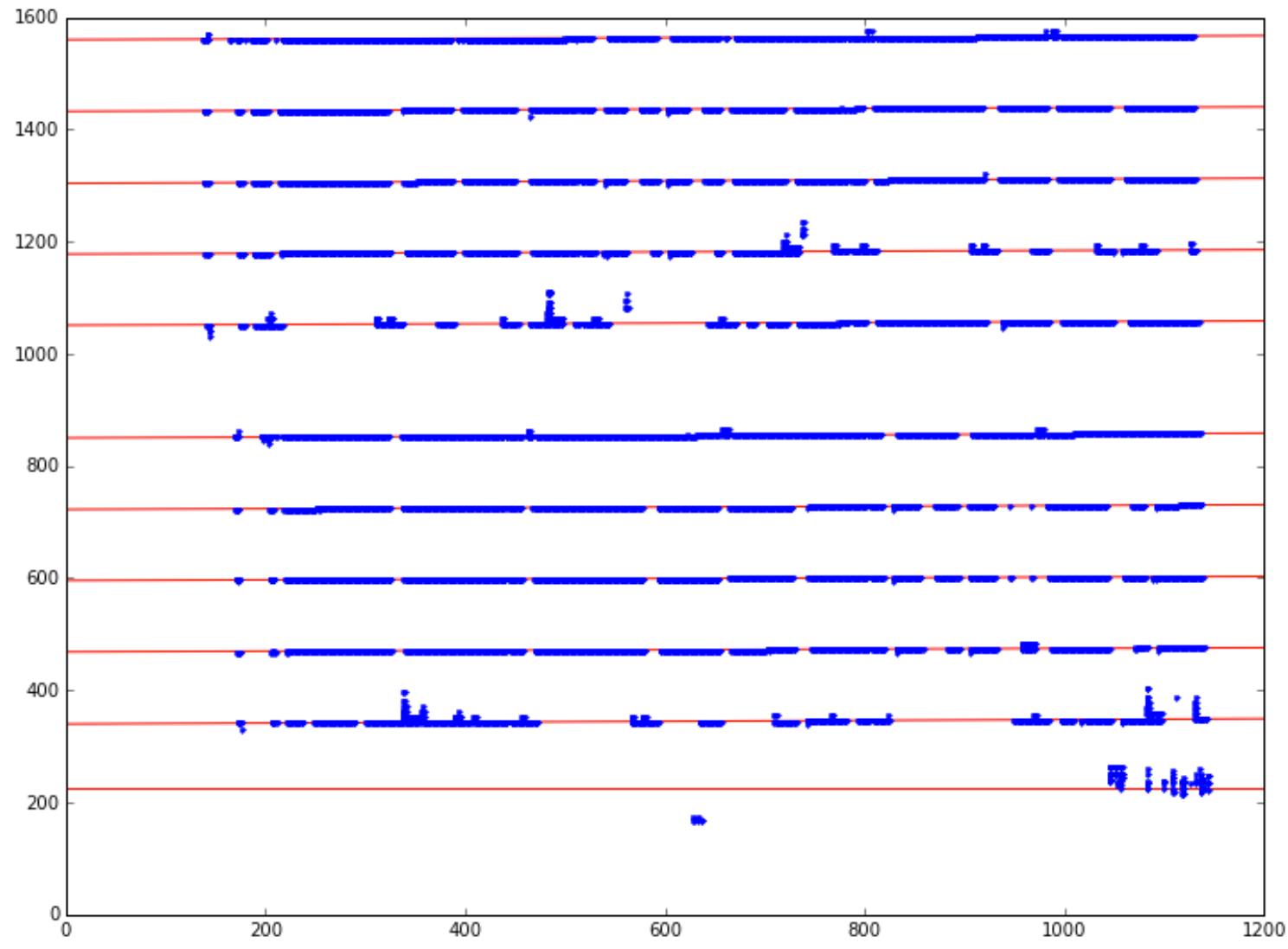
- S.** Soprano: The lyrics are "真的真的好想你 不是 假的假的好想你 好想 你 好想 你 好想". Red lines highlight notes in measures 1-3.
- A.** Alto: The lyrics are "oh ba ba la du du". Red lines highlight notes in measures 1-3.
- T.** Tenor: The lyrics are "du da du du". Red lines highlight notes in measures 1-3.
- T.** Tenor: The lyrics are "du du". Red lines highlight notes in measures 4-5.
- B.** Bass: The lyrics are "du du". Red lines highlight notes in measures 4-5.



Multiple Lines Fitting

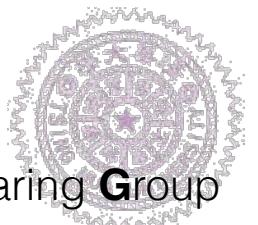


Multiple Lines Fitting



Multiple Lines Fitting

- Dynamic Programming: $O((\frac{N}{W})^2 \times W)$
- Hough Transform: $O(L \times N)$
- **RANSAC**: $O(N_{\text{iter}} \times N)$

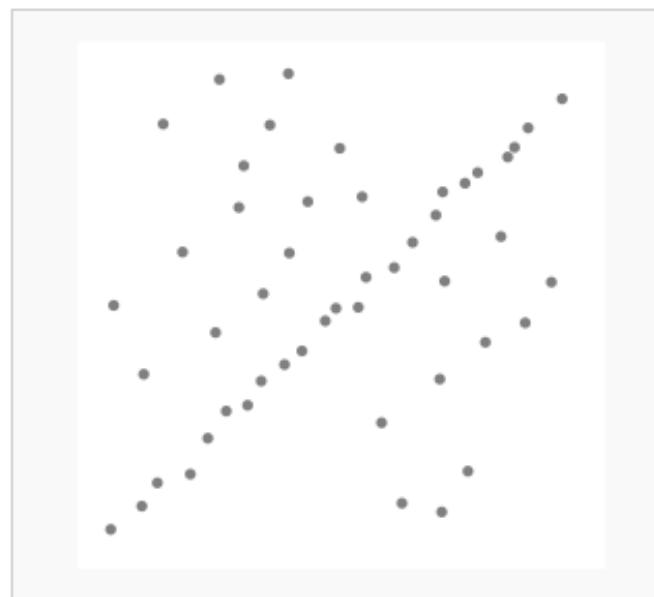


Multiple Lines Fitting

- RANSAC (RANdom SAmple Consensus)
 - single model fitting method
 - robust to outliers

M.A. Fischler and R.C. Bolles. **Random sample consensus: A paradigm for model fitting with applications to image analysis and automated cartography.**

Communications of the ACM, 24(6):381–395, 1981.



A data set with many outliers for which a line has to be fitted.

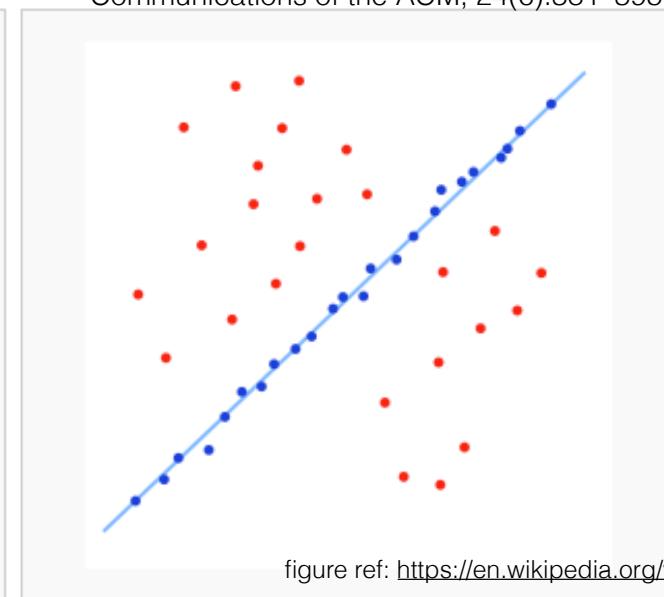
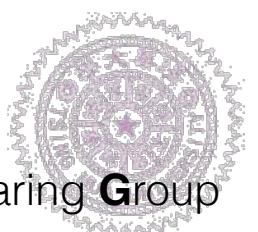
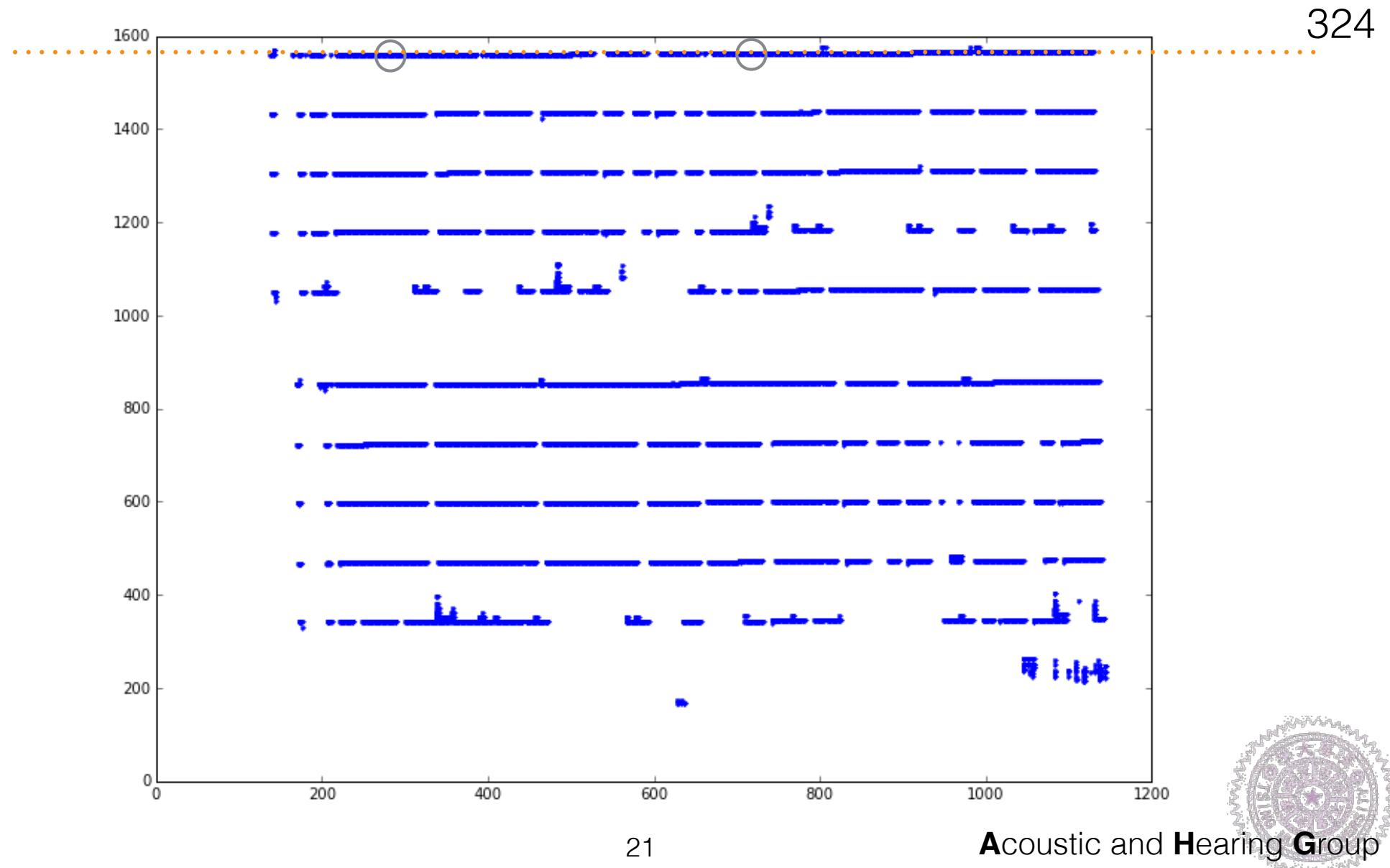


figure ref: <https://en.wikipedia.org/wiki/RANSAC>

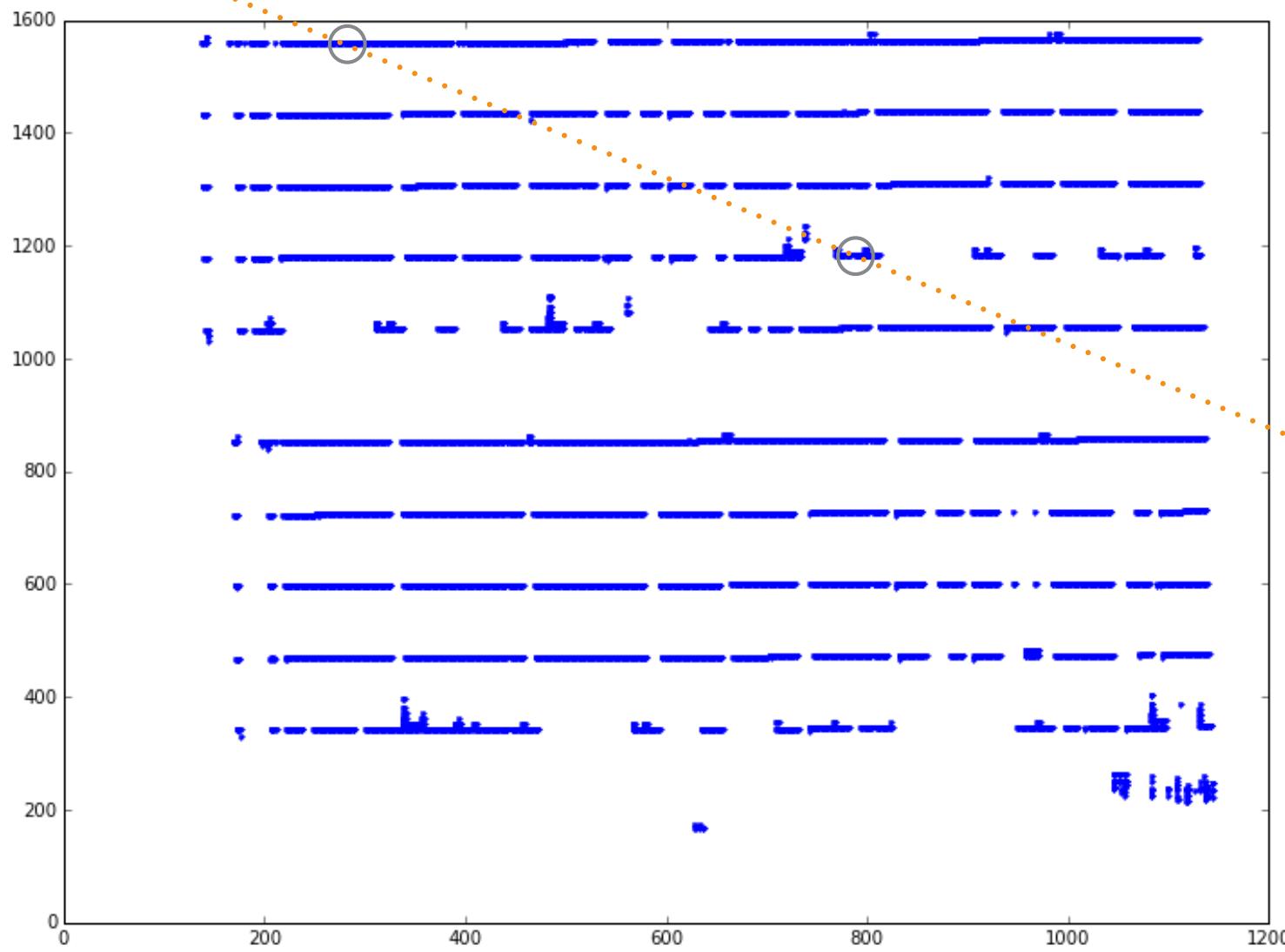
Fitted line with RANSAC; outliers have no influence on the result.



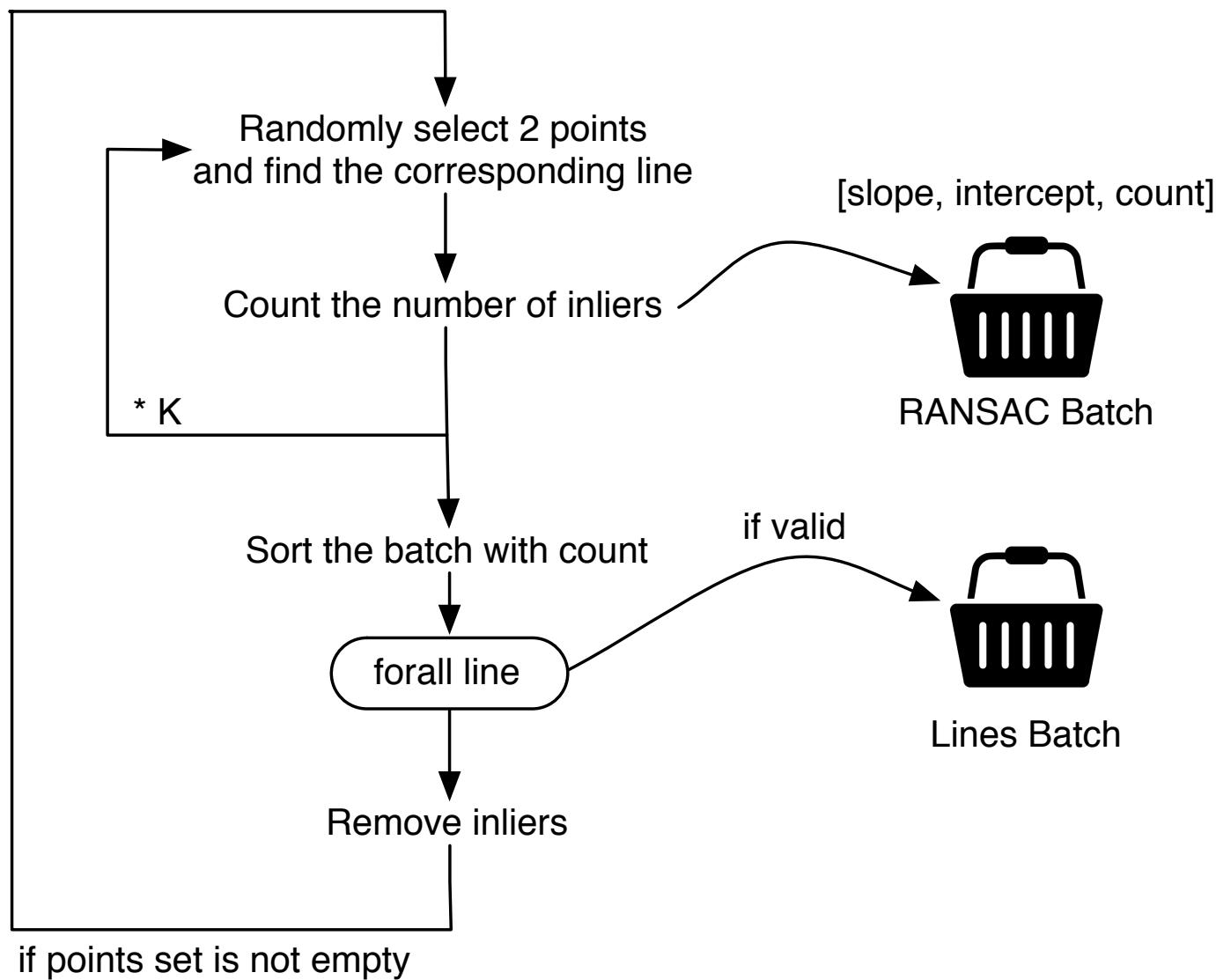
RANSAC: Example



RANSAC: Example



Iterative RANSAC



Detected Staves

40 7

S. 說來說去都只想讓你開心 好想你 好想你 好想你 好想你 是

A. 說來說去都只想讓你開心 du du

T. 說來說去都只想讓你開心 du

T. 00 00 00 du

B. dm dm



Staff Images (Before Staff Removal)

$$l_i = c_y + (i - 3) \times (w + s)$$

S. 40 說來說去都只想讓 你開心 好相你 好相你 好相你 好相你 早

A. 說來說去都只想讓 你開心 du du

T. 說來說去都只想讓 你開心 du

T. oo oo oo du

Staff Images (After Staff Removal)

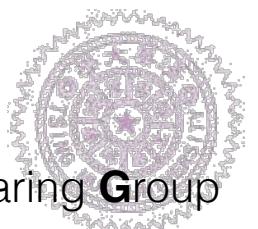
$$l_i = c_y + (i - 3) \times (w + s)$$

The image displays four staves of musical notation, likely for a vocal performance. Each staff is represented by a vertical line with a clef and key signature, followed by a series of notes and rests. The notes are black with stems, and the rests are white with diagonal lines. Below each staff, there is Chinese lyrics and some accompanying text. The staves are arranged vertically, separated by horizontal lines.

- Soprano (S.)**: The first staff starts with a treble clef and a key signature of one flat. It contains lyrics: "說來說去都只想讓 你開心 好相 你好相 你好相 你好相 你好相 你好相 是你".
- Alto (A.)**: The second staff starts with a treble clef and a key signature of one flat. It contains lyrics: "說來說去都只想讓 你開心 du du".
- Tenor (T.)**: The third staff starts with a treble clef and a key signature of one flat. It contains lyrics: "說來說去都只想讓 你開心 du".
- Bass (B.)**: The fourth staff starts with a bass clef and a key signature of one flat. It contains lyrics: "oo oo oo du".

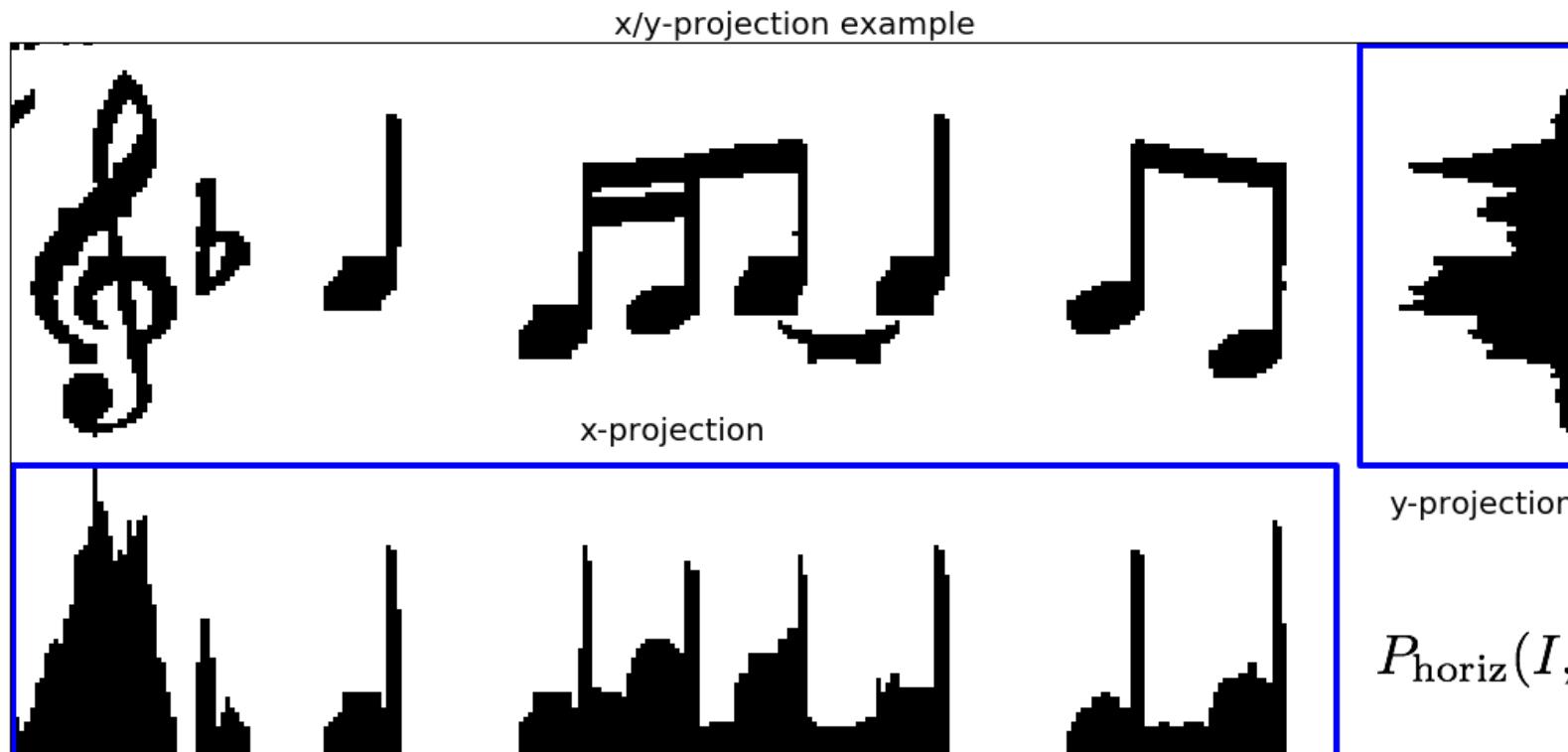


Recognition



Acoustic and Hearing Group

Projection



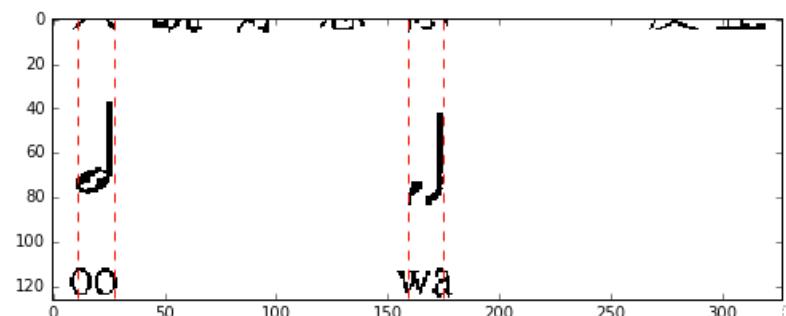
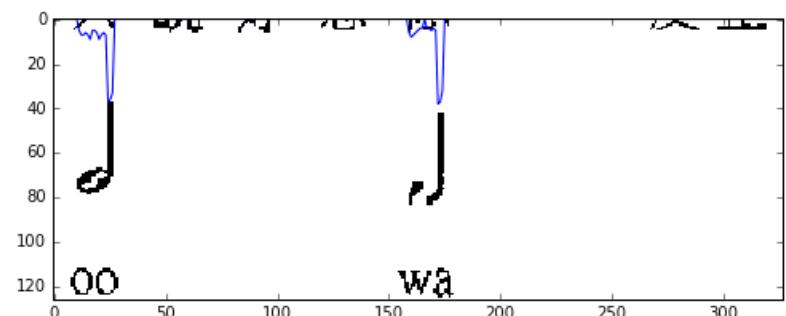
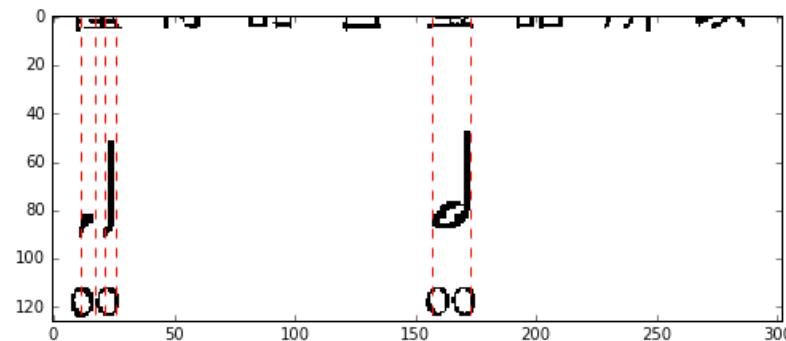
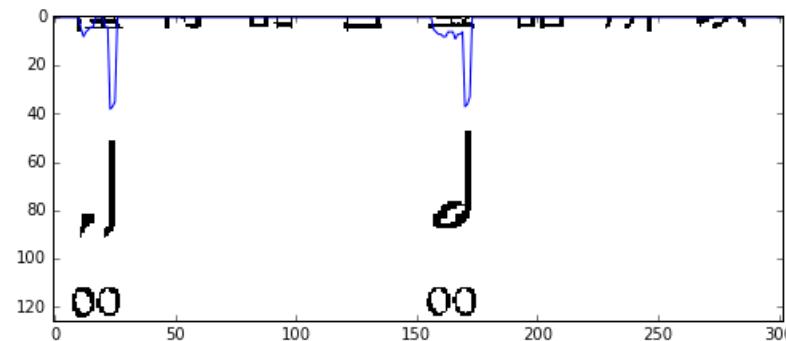
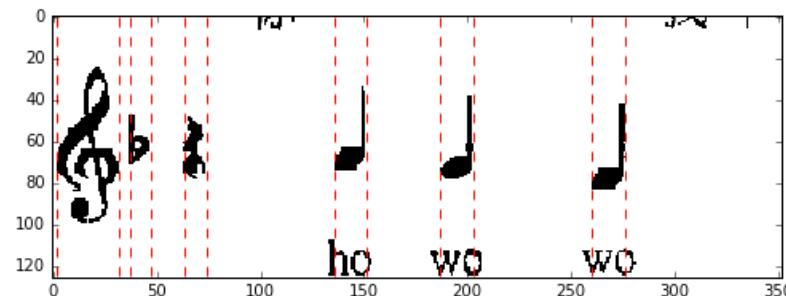
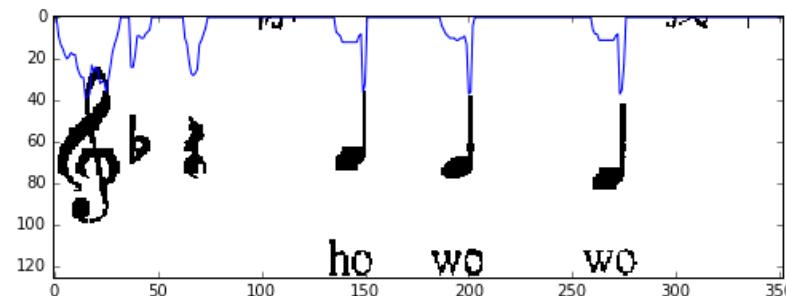
$$P_{\text{horiz}}(I, x_0) = \sum_y I(x_0, y)$$

$$P_{\text{vert}}(I, y_0) = \sum_x I(x, y_0)$$



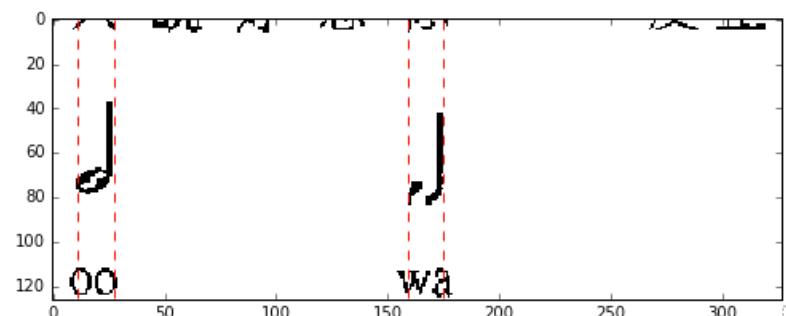
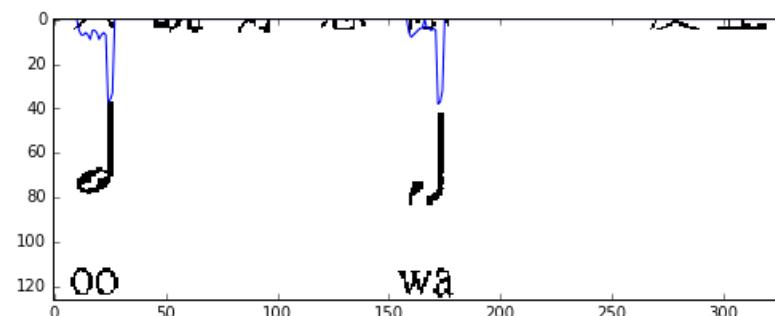
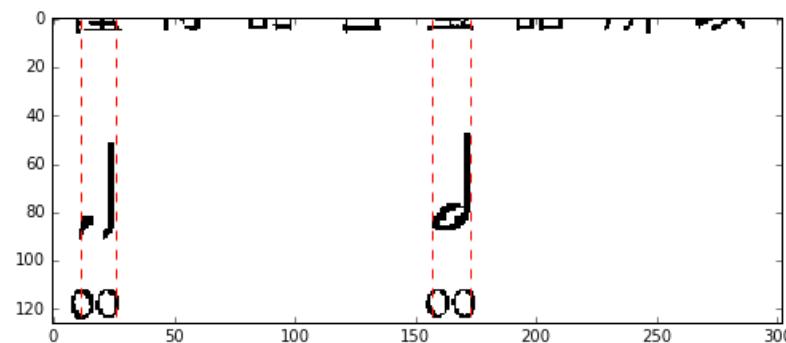
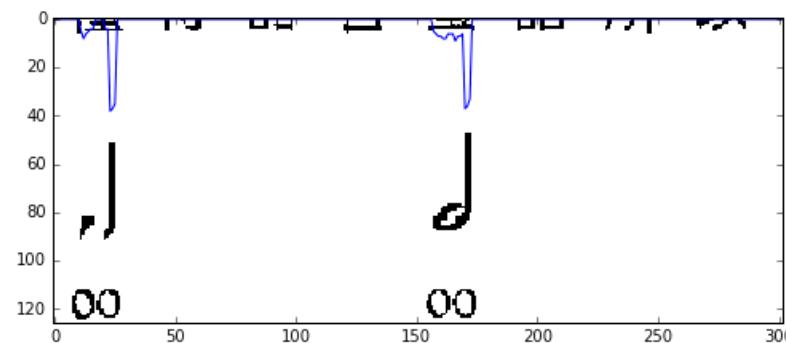
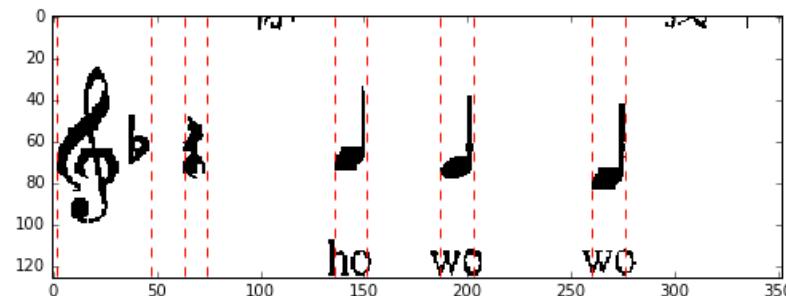
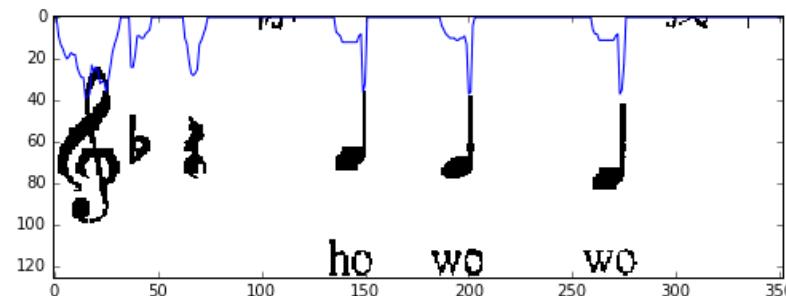
Onset Detection

merge_thresh=0



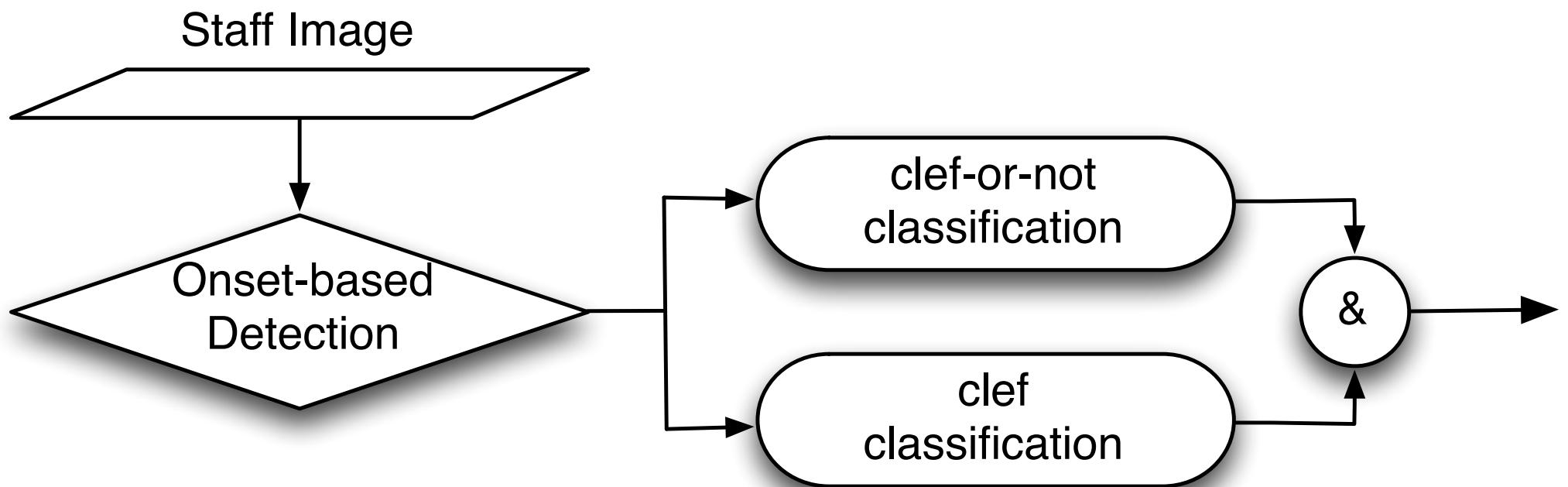
Onset Detection

merge_thresh=10

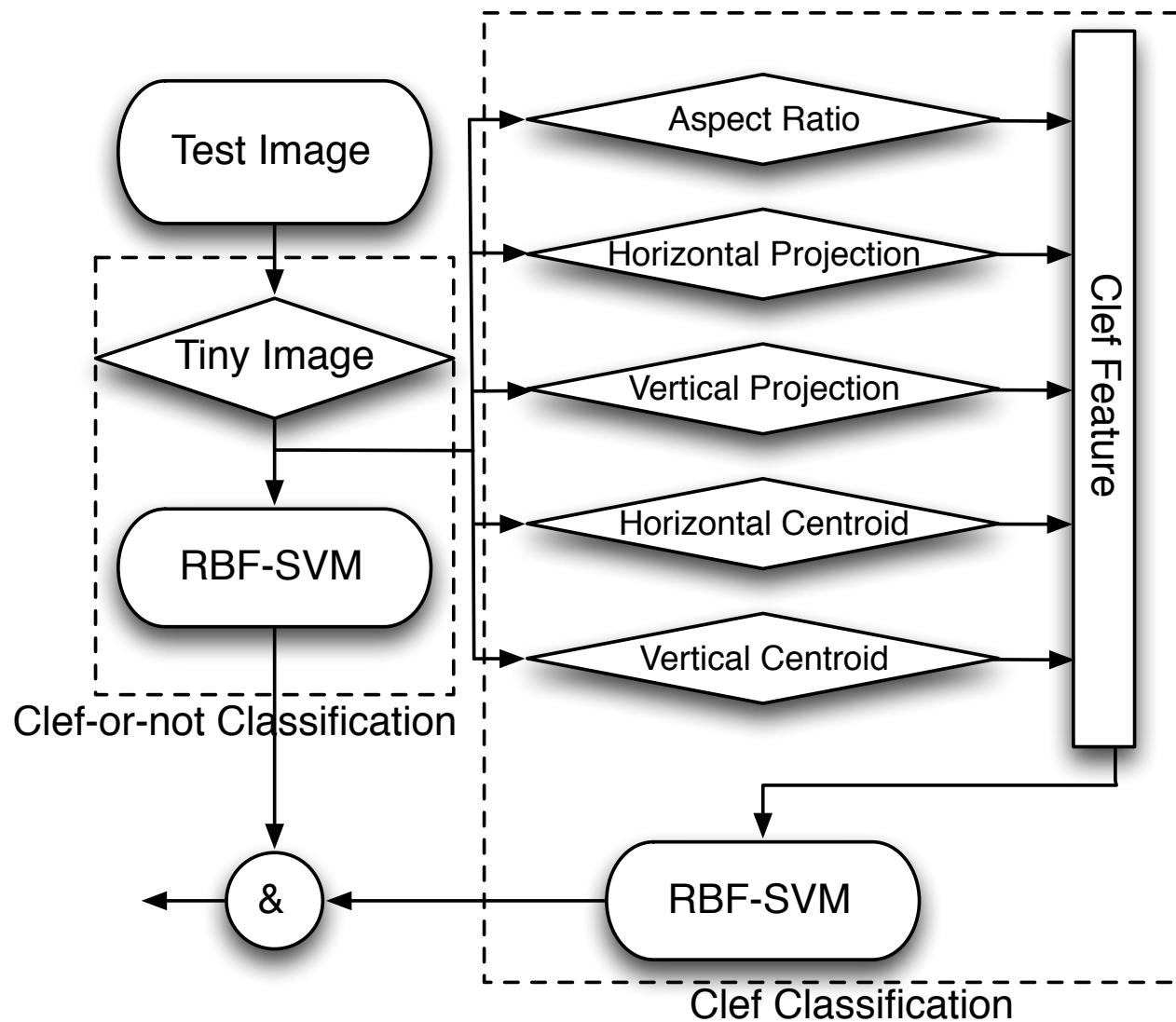


Recognition Pitch-unrelated Symbols

Clef Detection



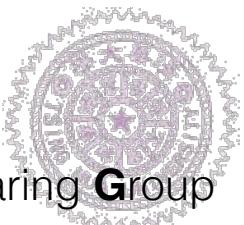
Clef Detection: Feature



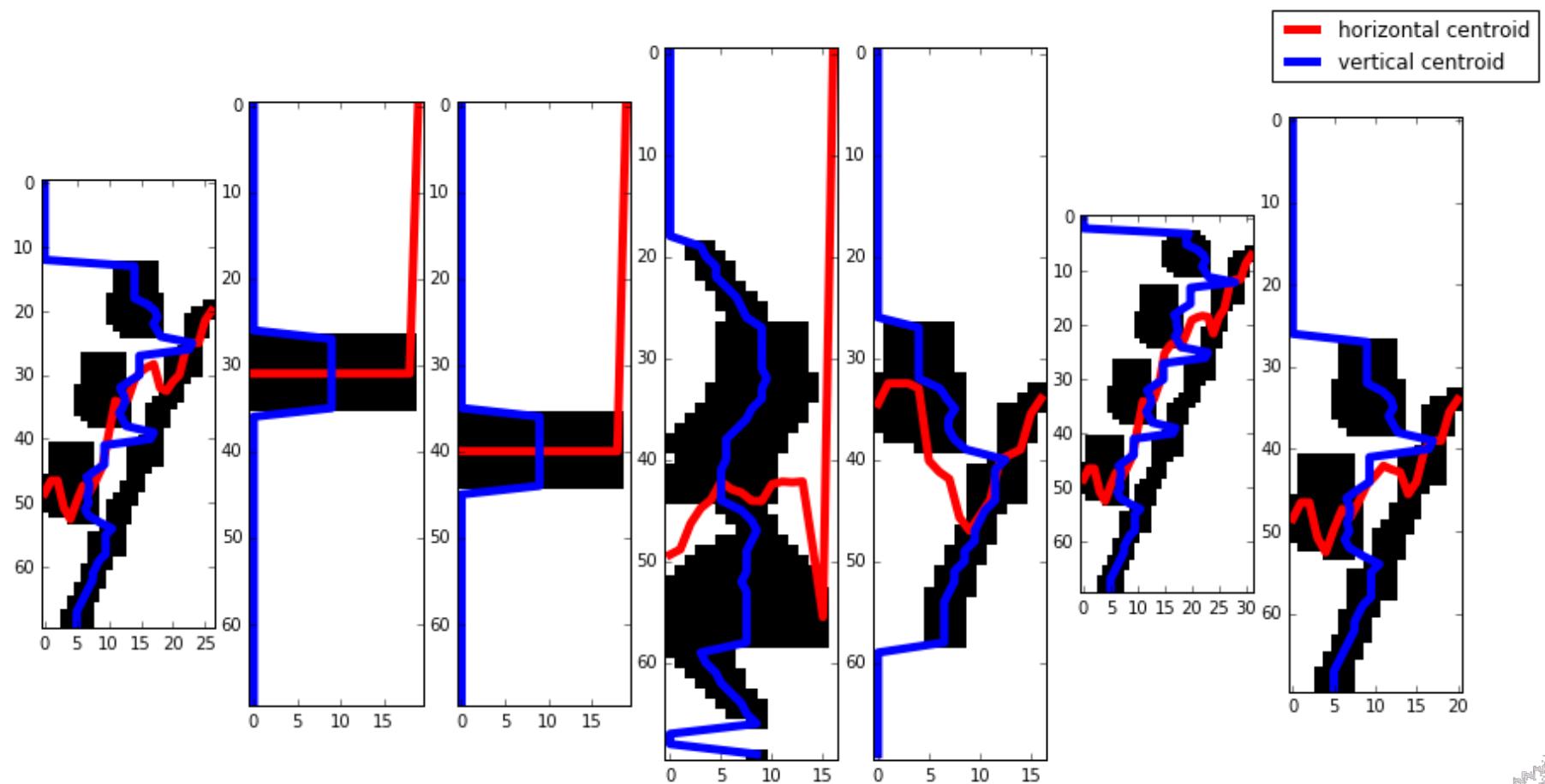
Centroid

$$C_{\text{horiz}}(x_0) = \frac{\sum_y I(x_0, y) \times y}{\sum_y I(x_0, y)}$$

$$C_{\text{vert}}(y_0) = \frac{\sum_x I(x, y_0) \times x}{\sum_x I(x, y_0)}$$



Centroid

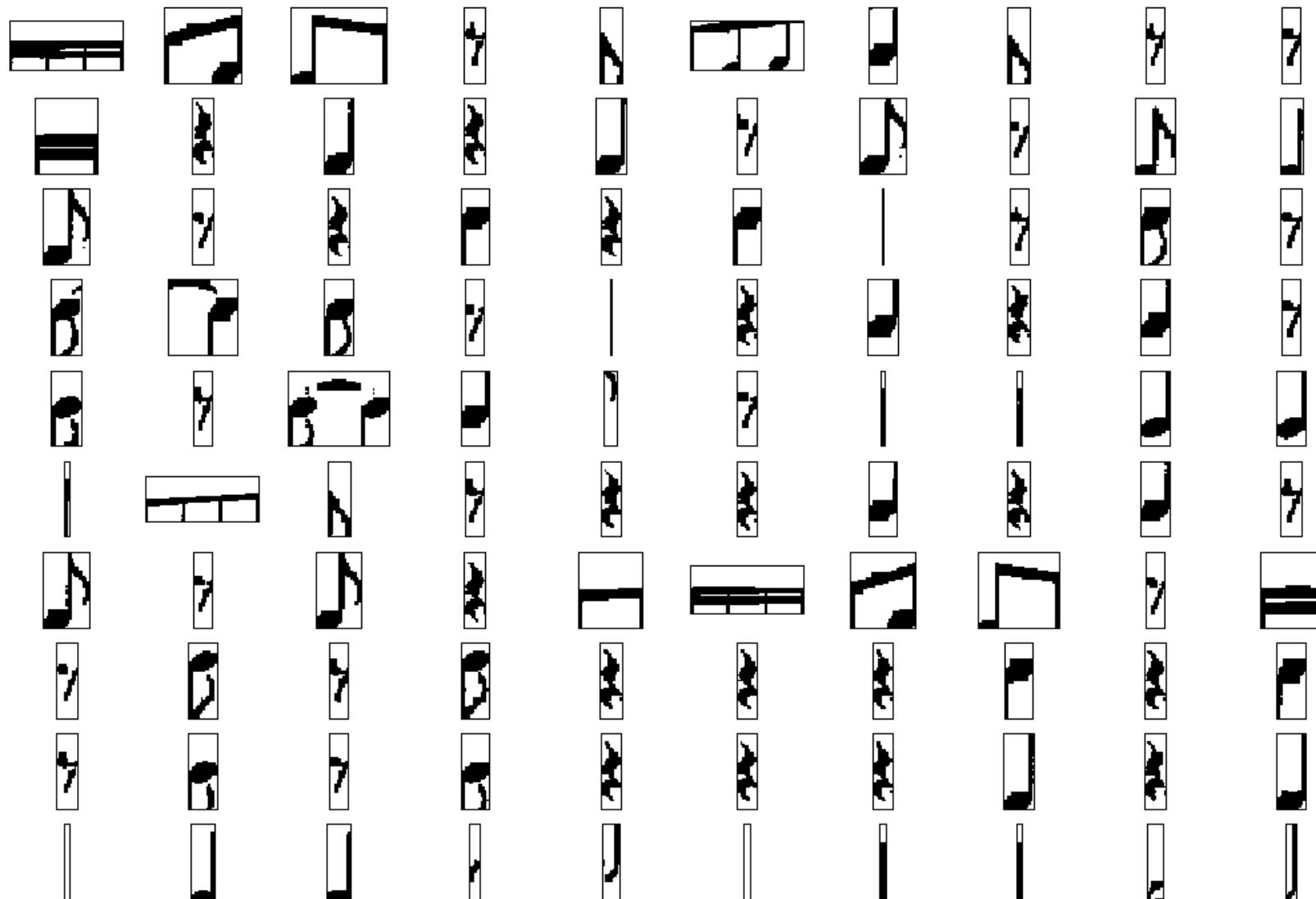


$$C_{\text{horiz}}(x_0) = \frac{\sum_y I(x_0, y) \times y}{\sum_y I(x_0, y)}$$

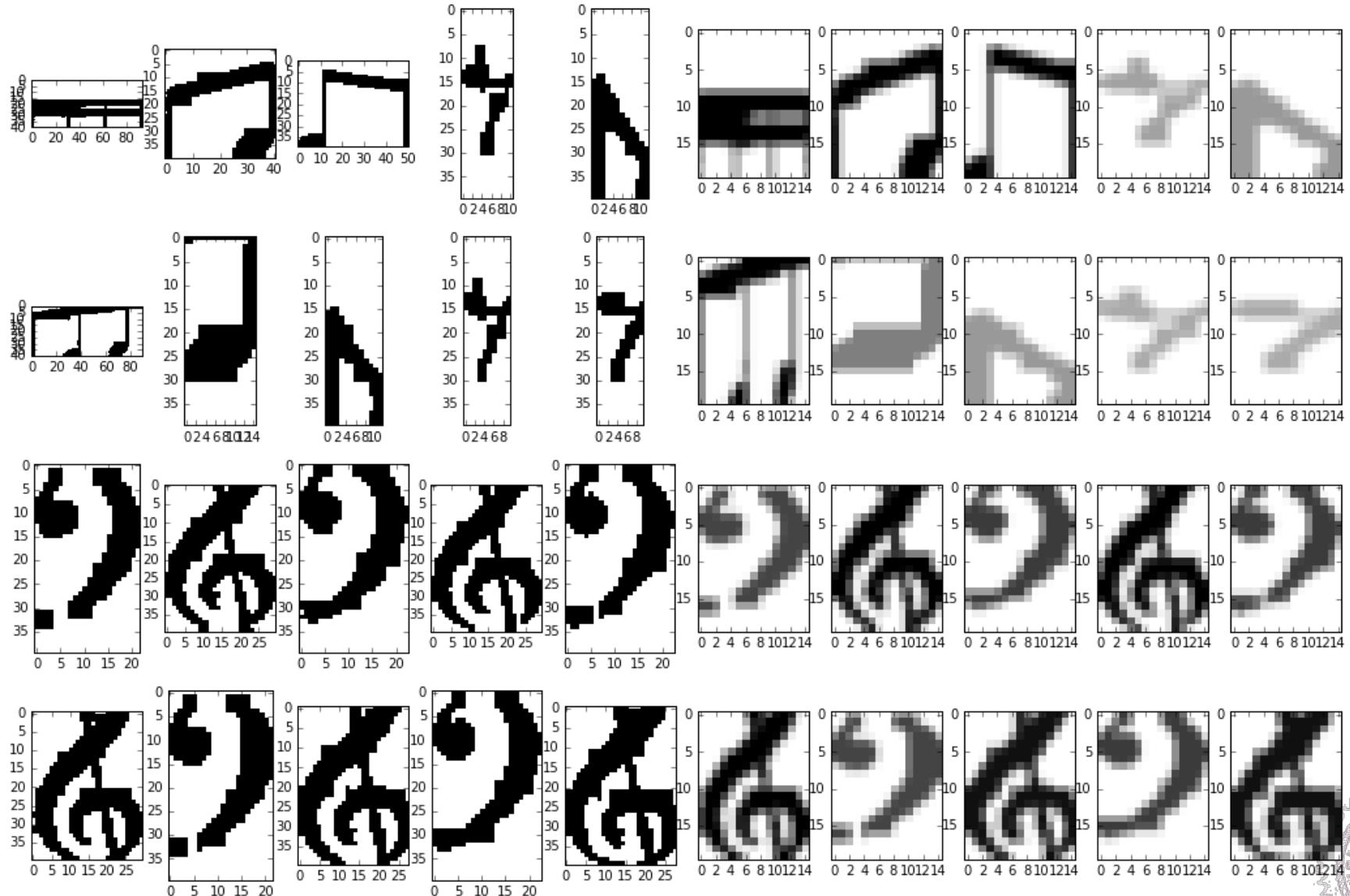
$$C_{\text{vert}}(y_0) = \frac{\sum_x I(x, y_0) \times x}{\sum_x I(x, y_0)}$$



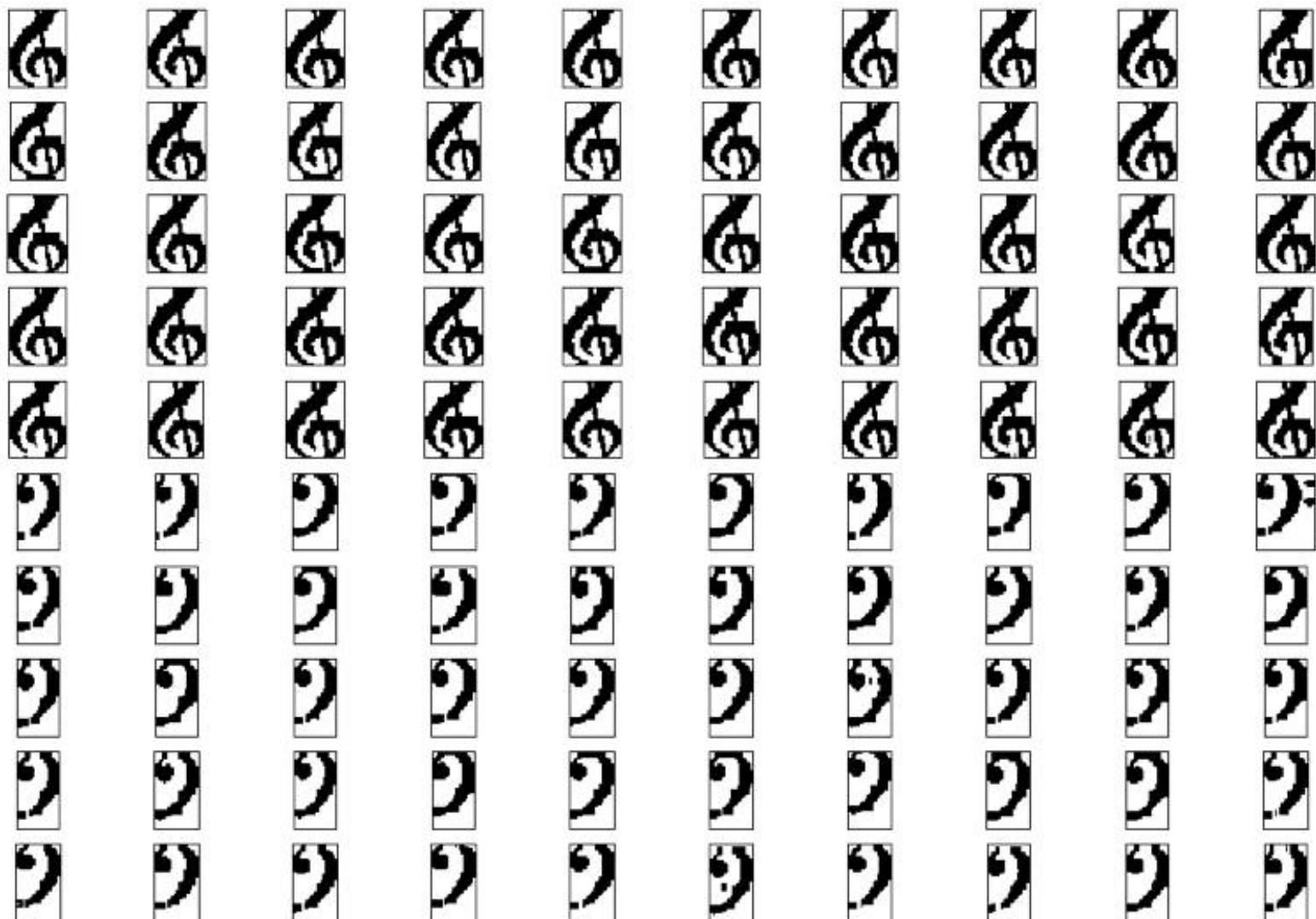
Clef-or-not Samples



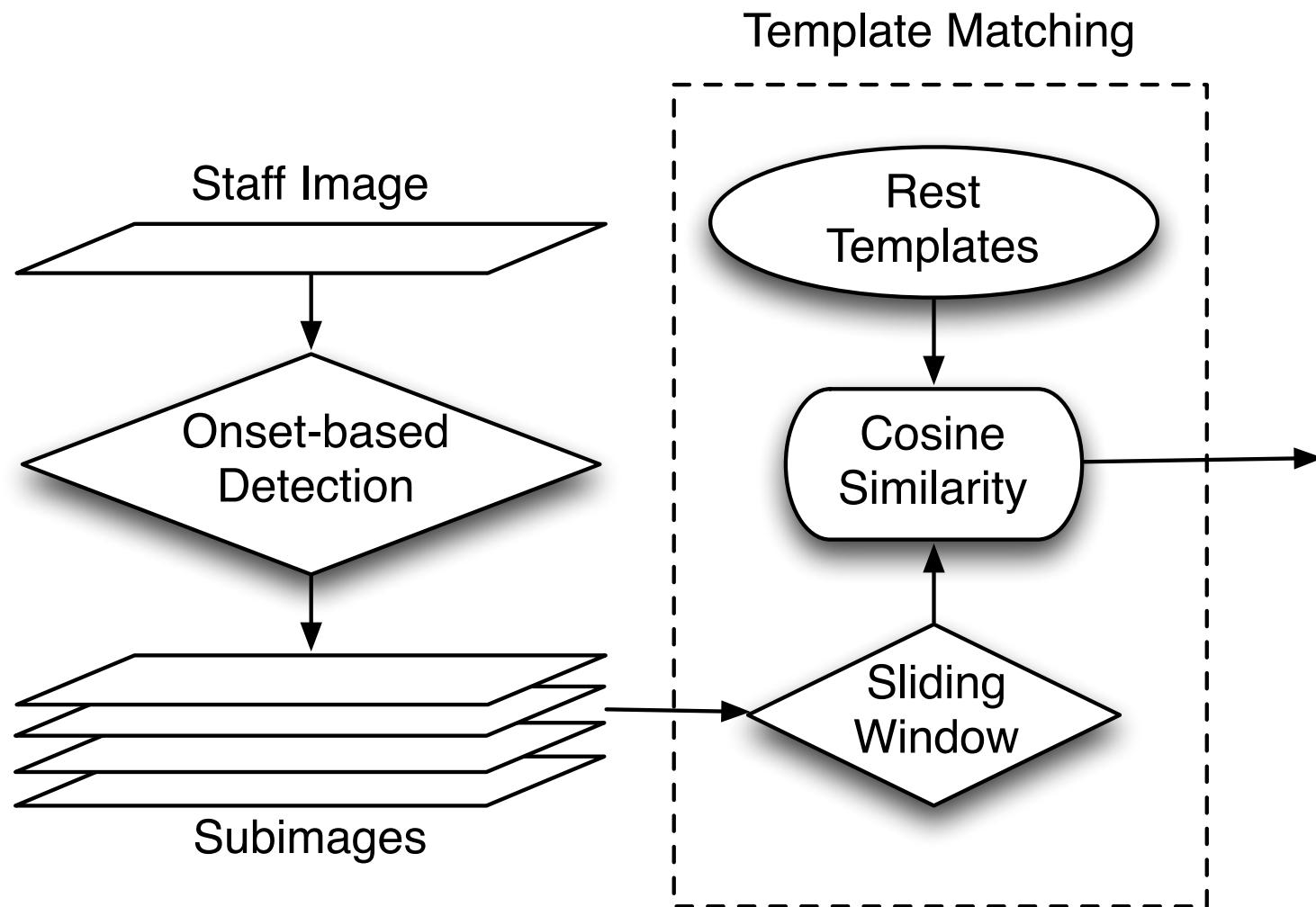
Clef-or-not Feature



Clef Samples

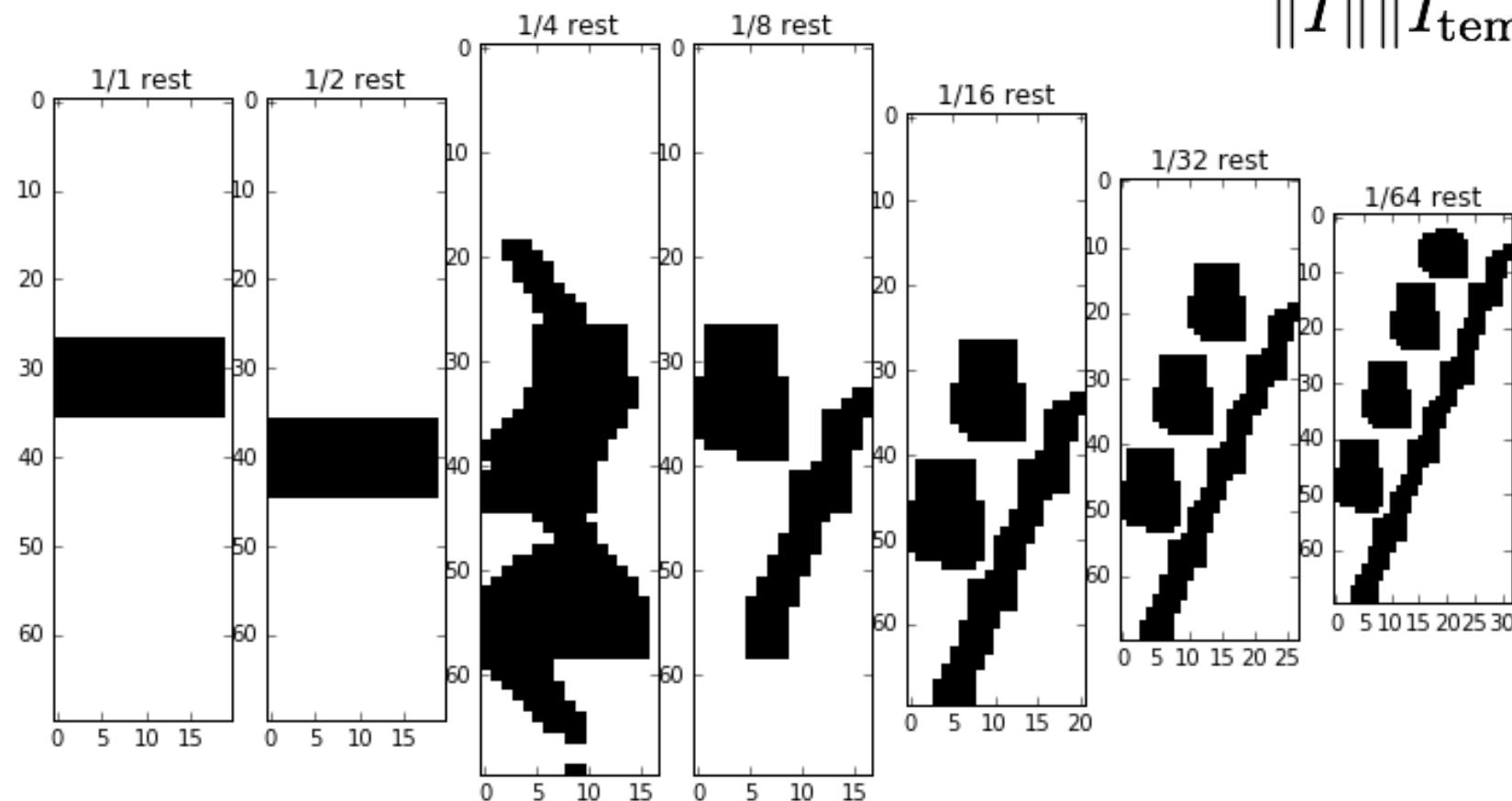


Rest Detection



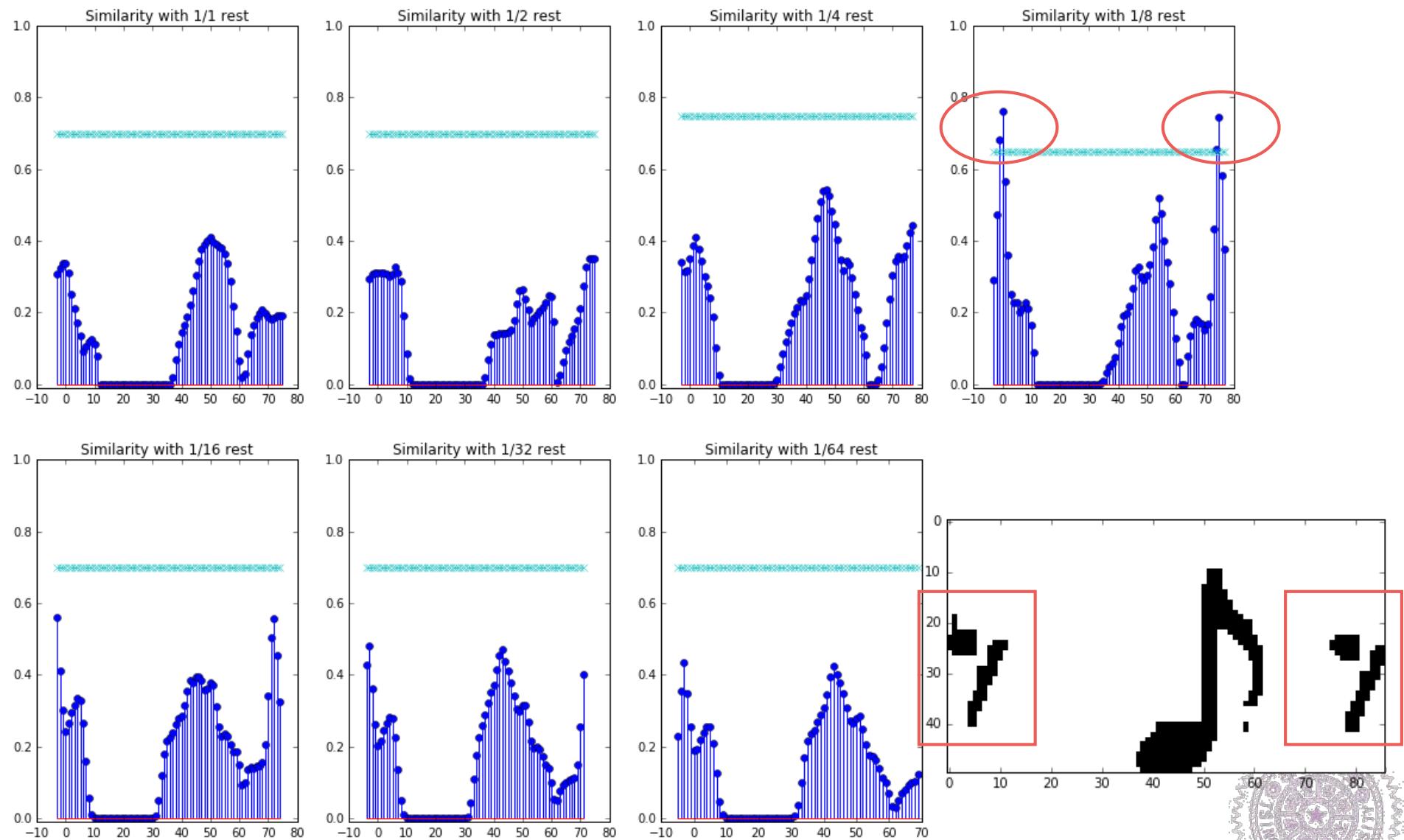
Rest Templates

$$\text{score} = \frac{\langle I, I_{\text{template}} \rangle}{\|I\| \|I_{\text{template}}\|}$$



Rest Detection

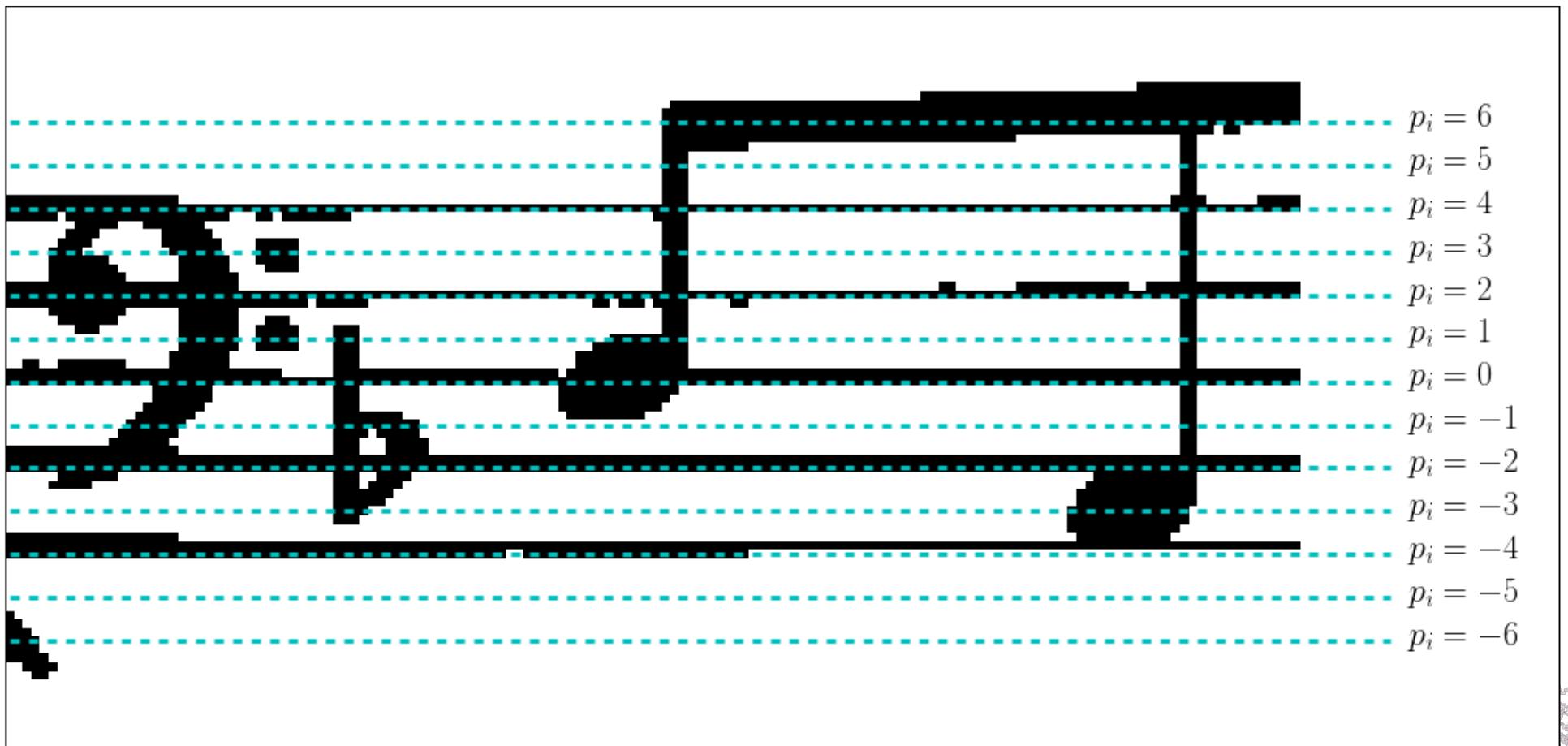
Template Matching



Recognition Pitch-related Symbols

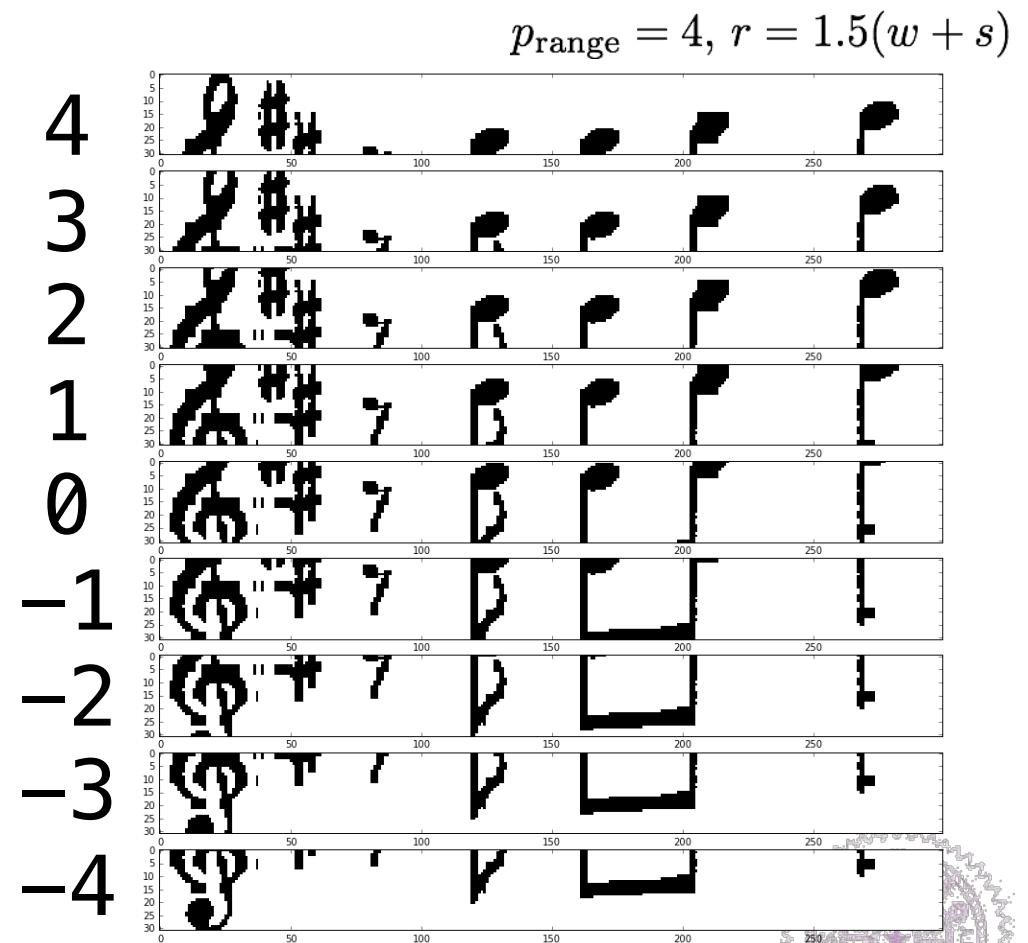
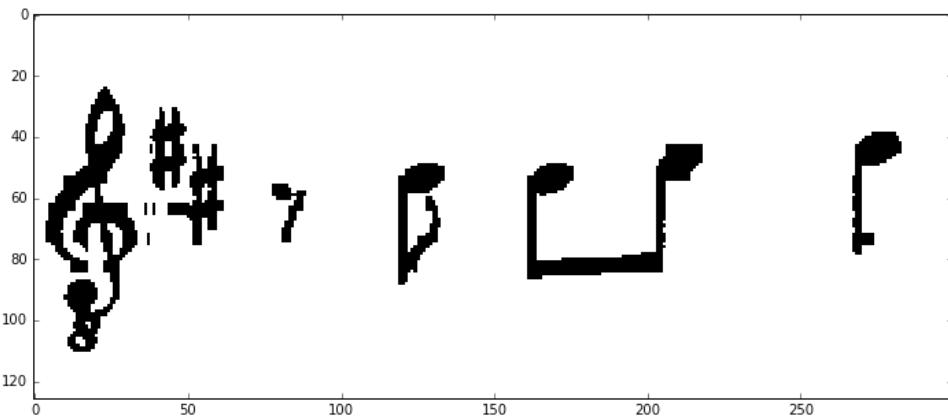
Pitch Cropping

$$y_{p_i} = l_3 - p_i \times 0.5(w + s)$$

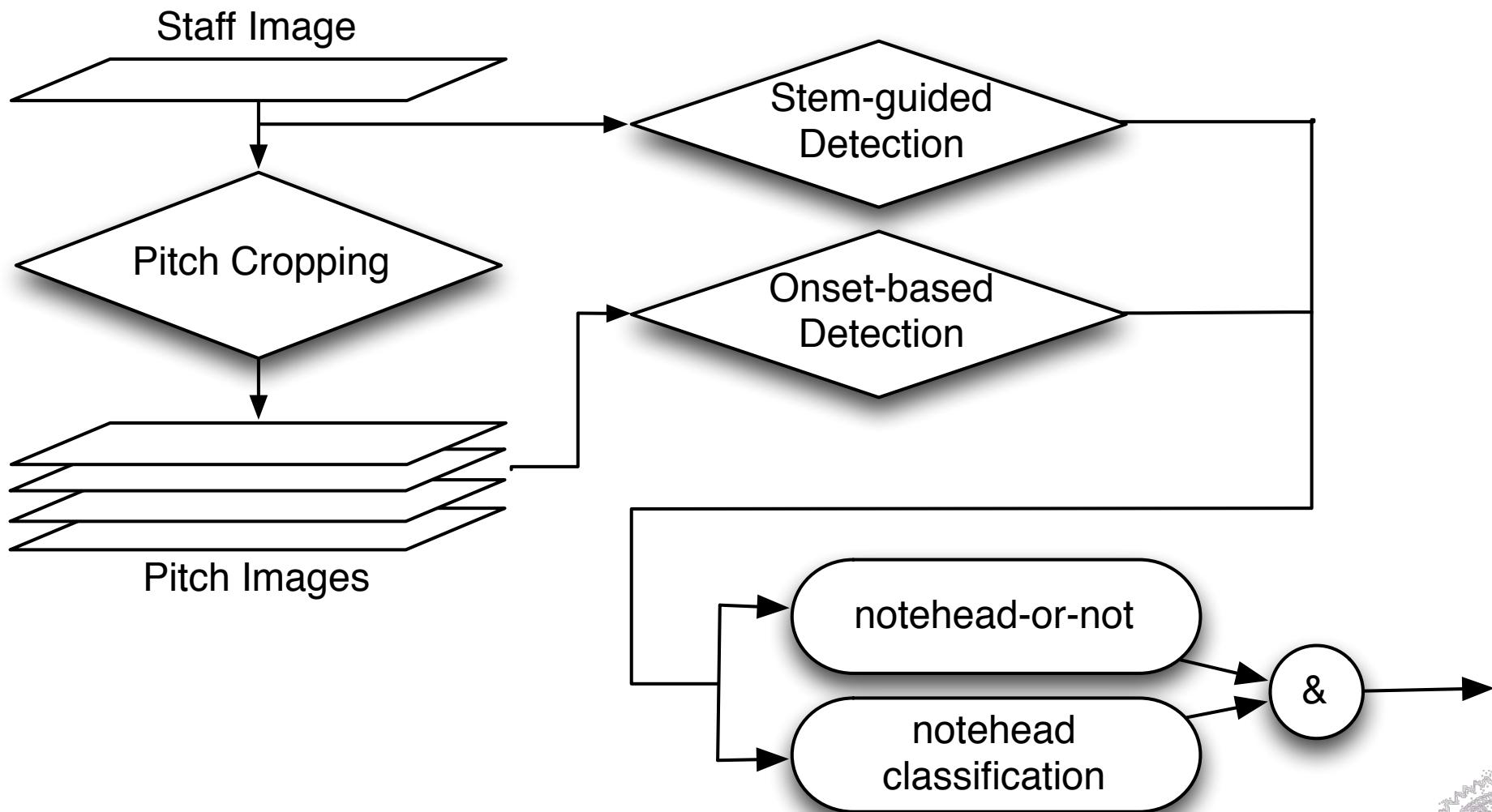


Pitch Cropping

$$\rho_{\text{seg}}(I, p_{\text{range}}, r) = \{I(:, y_{p_i} - r : y_{p_i} + r) \mid p_i = -p_{\text{range}}, \dots, p_{\text{range}}\}$$

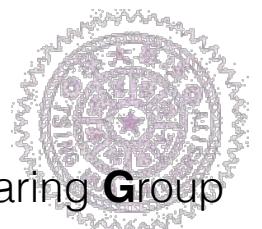
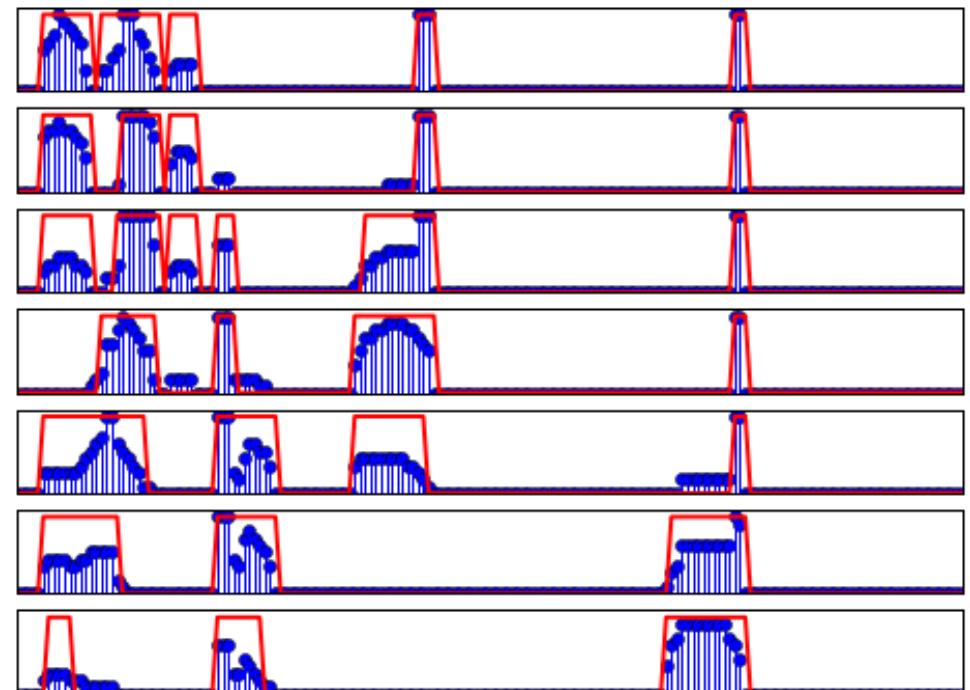
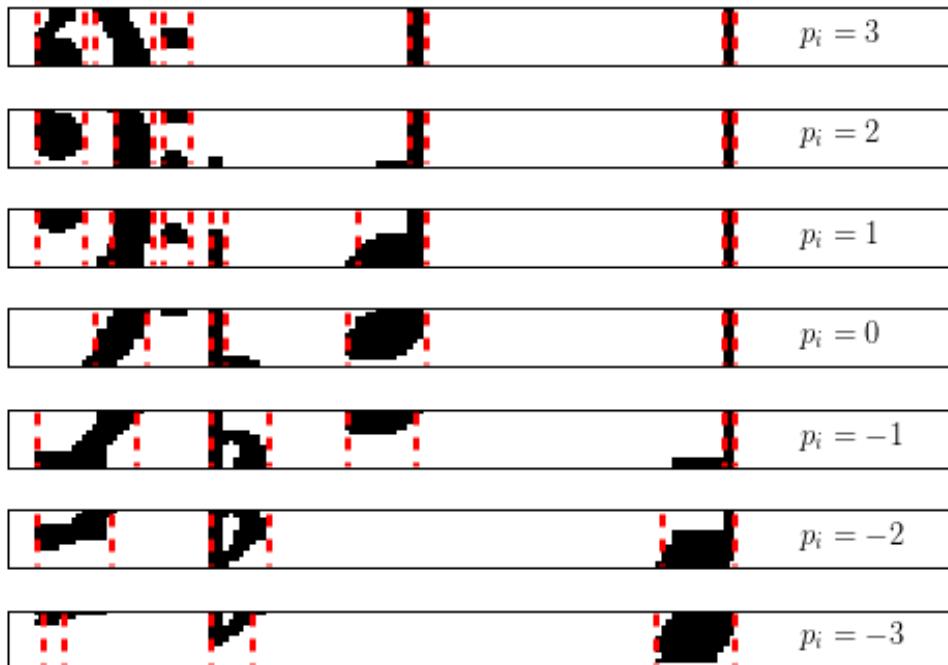


Notehead Detection



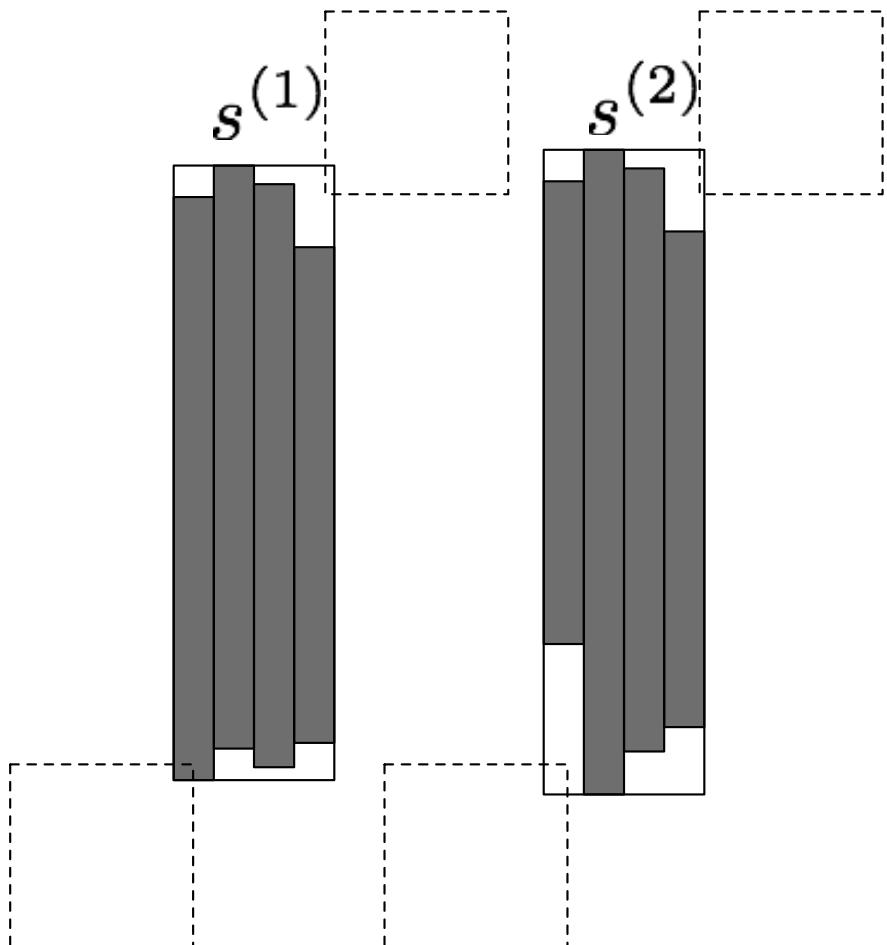
Notehead Detection

Onset-based



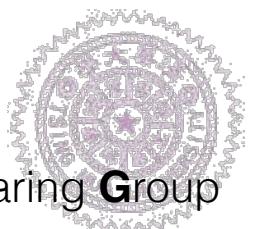
Notehead Detection

Stem-guided

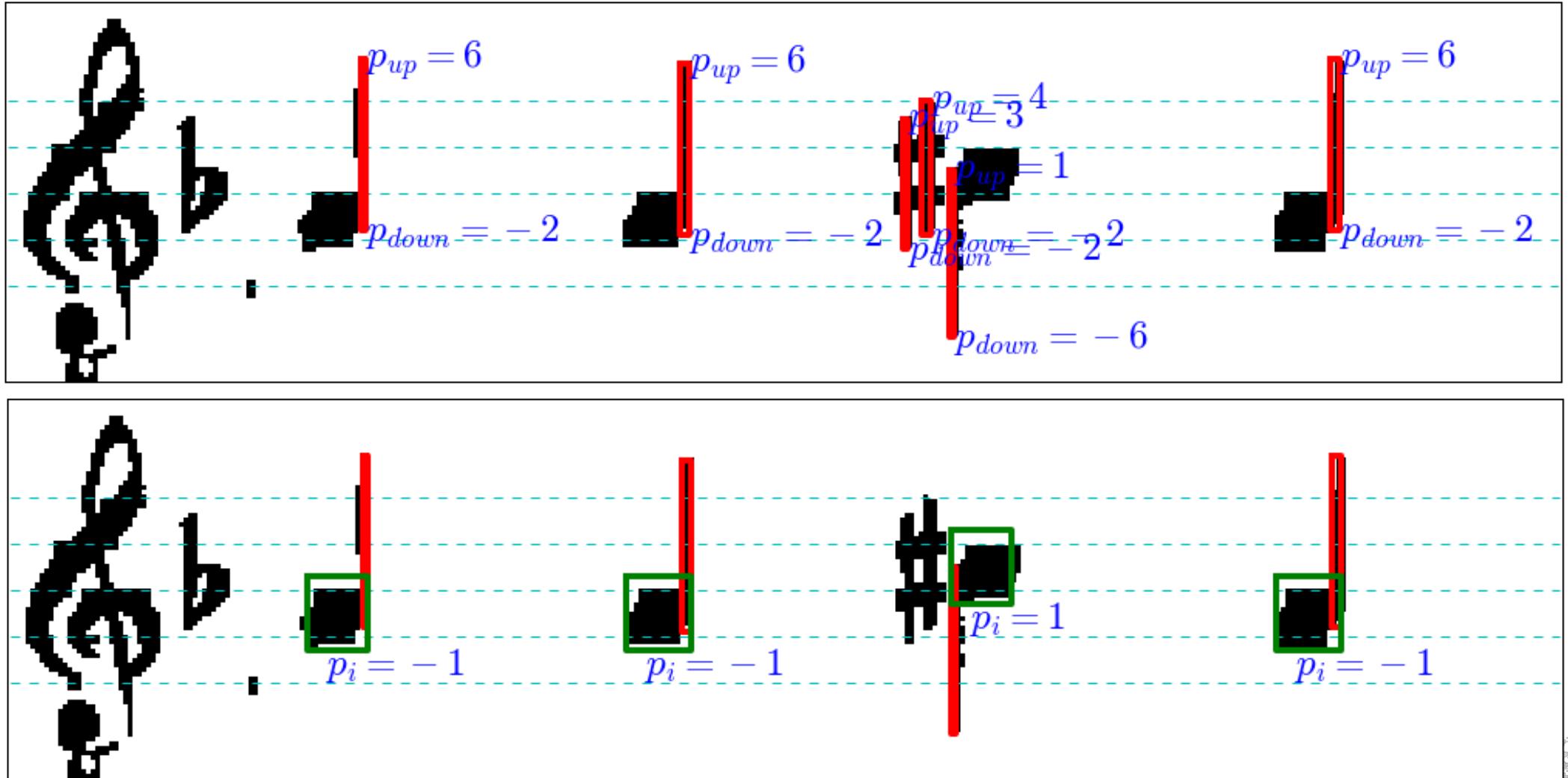


$$p_u = \text{round}\left(\frac{l_3 - y_{\min}}{0.5(w+s)}\right)$$

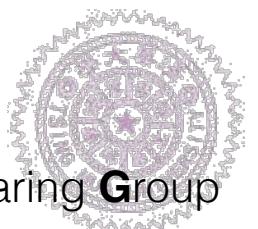
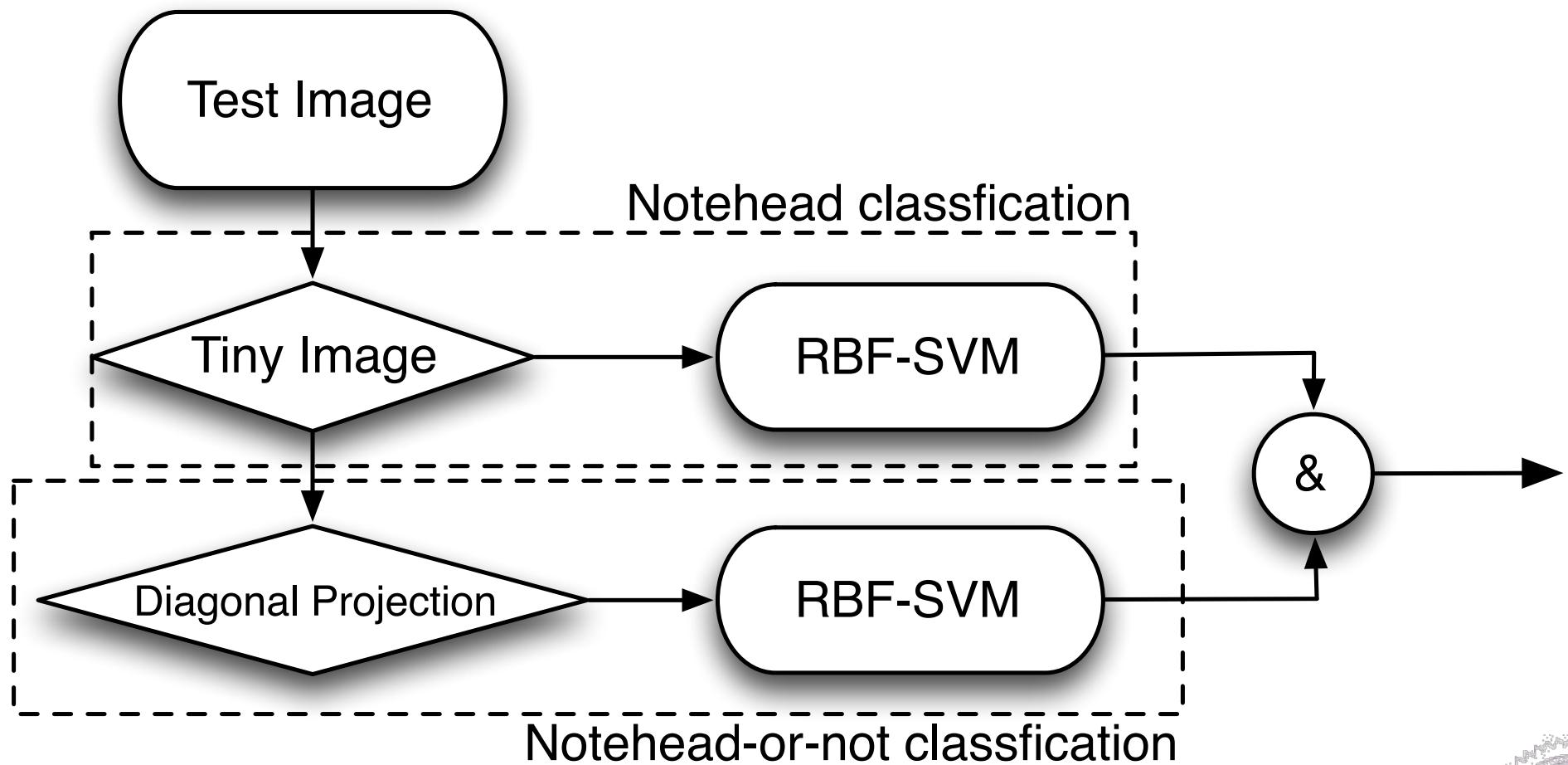
$$p_d = \text{round}\left(\frac{l_3 - y_{\max}}{0.5(w+s)}\right)$$



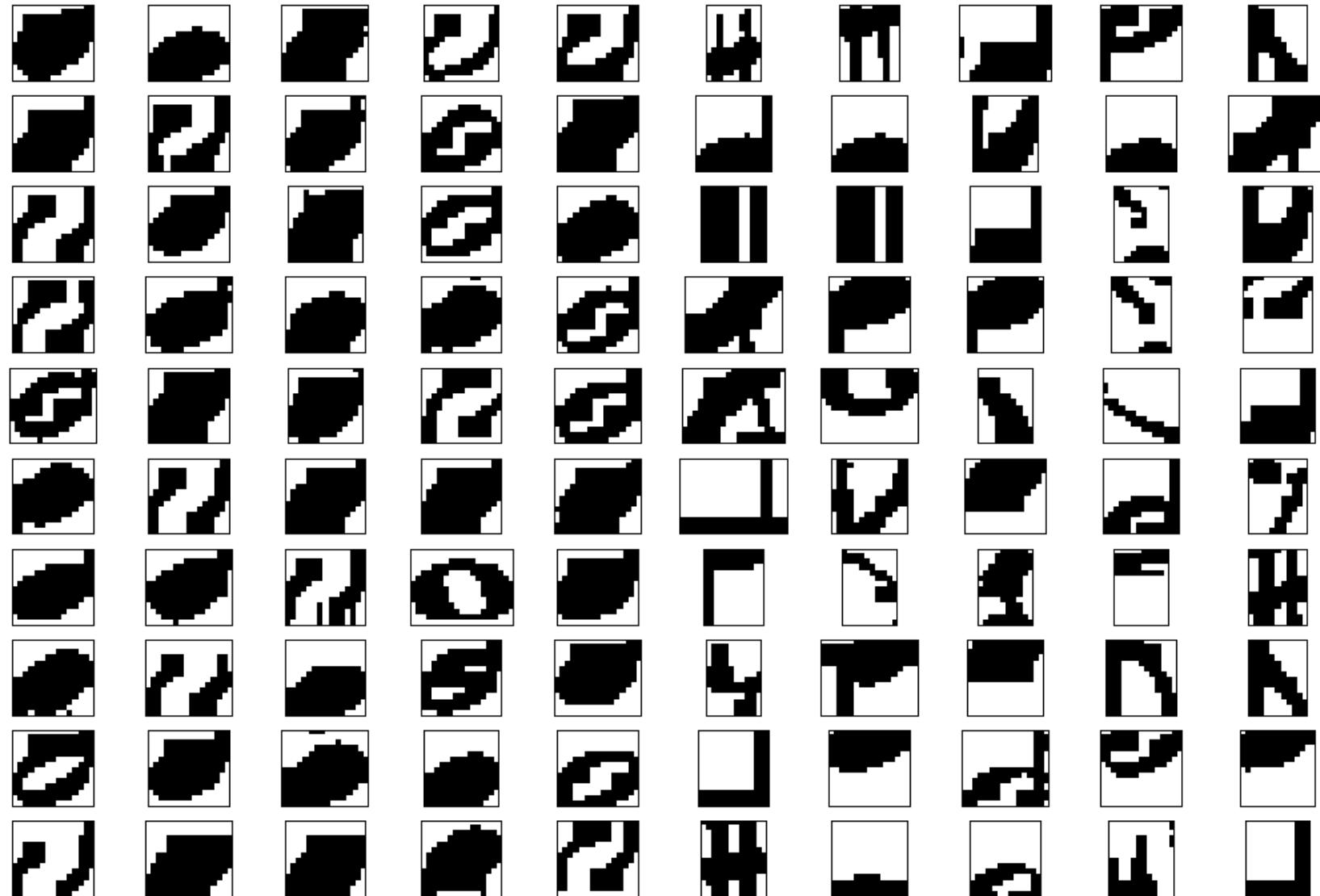
Notehead Detection Stem-guided



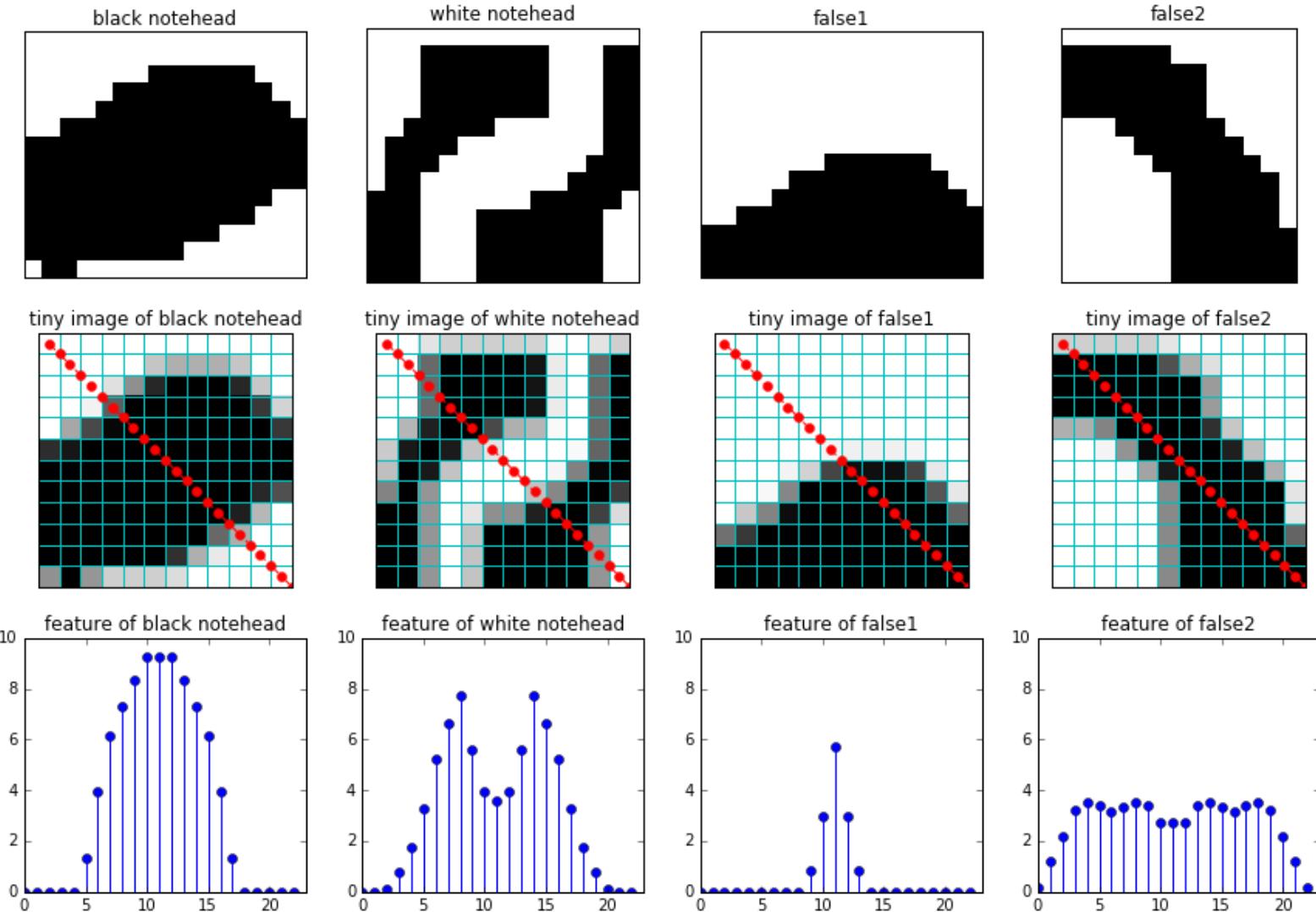
Notehead Detection: Feature



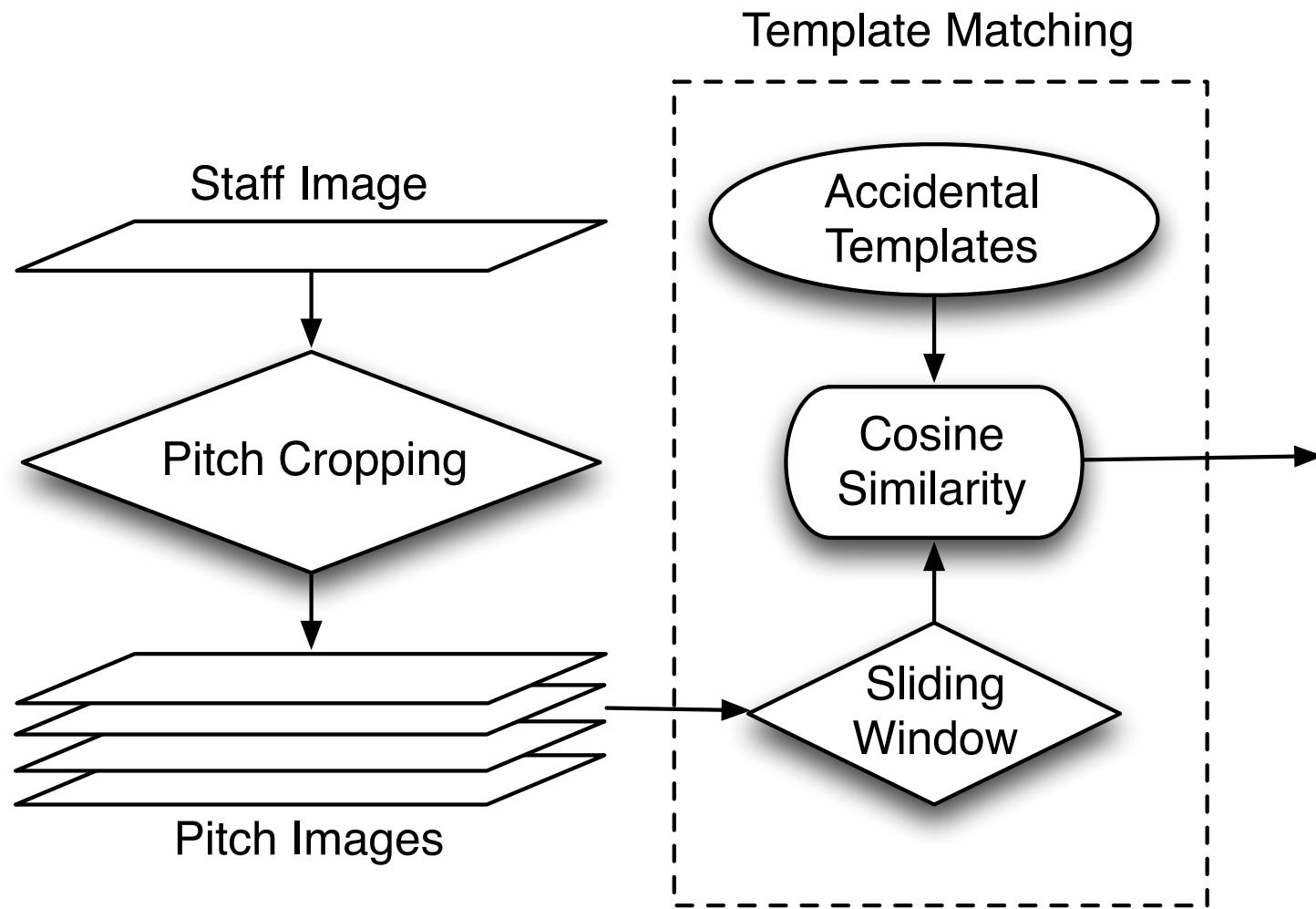
Notehead-or-not Samples



Notehead-or-not Feature

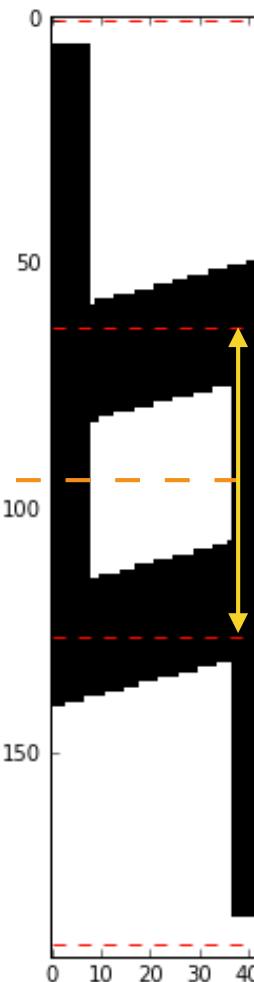
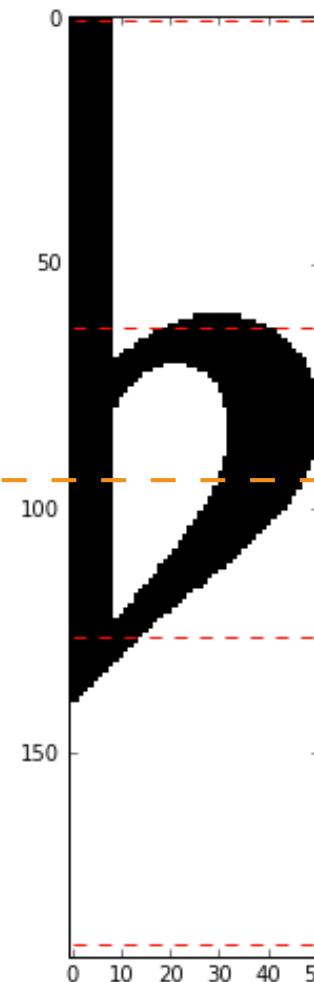
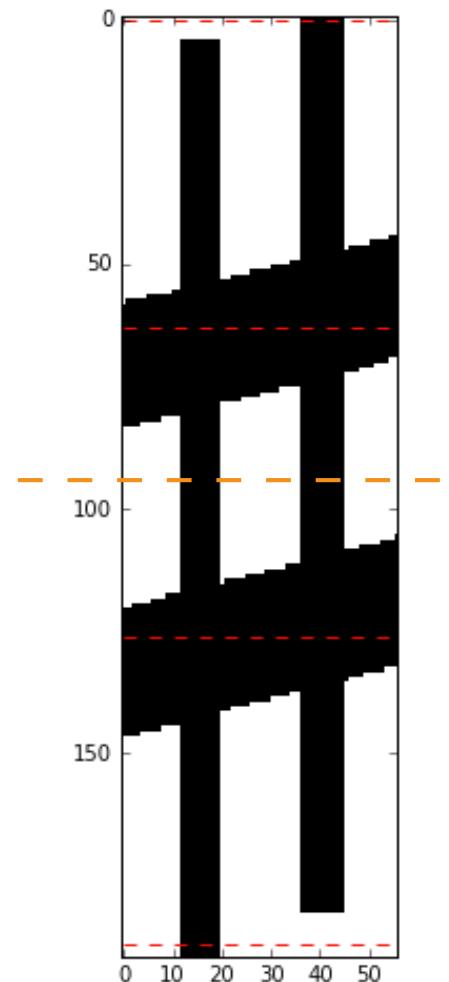


Accidental Detection



Accidental Templates

$$\text{score} = \frac{\langle I, I_{\text{template}} \rangle}{\|I\| \|I_{\text{template}}\|}$$

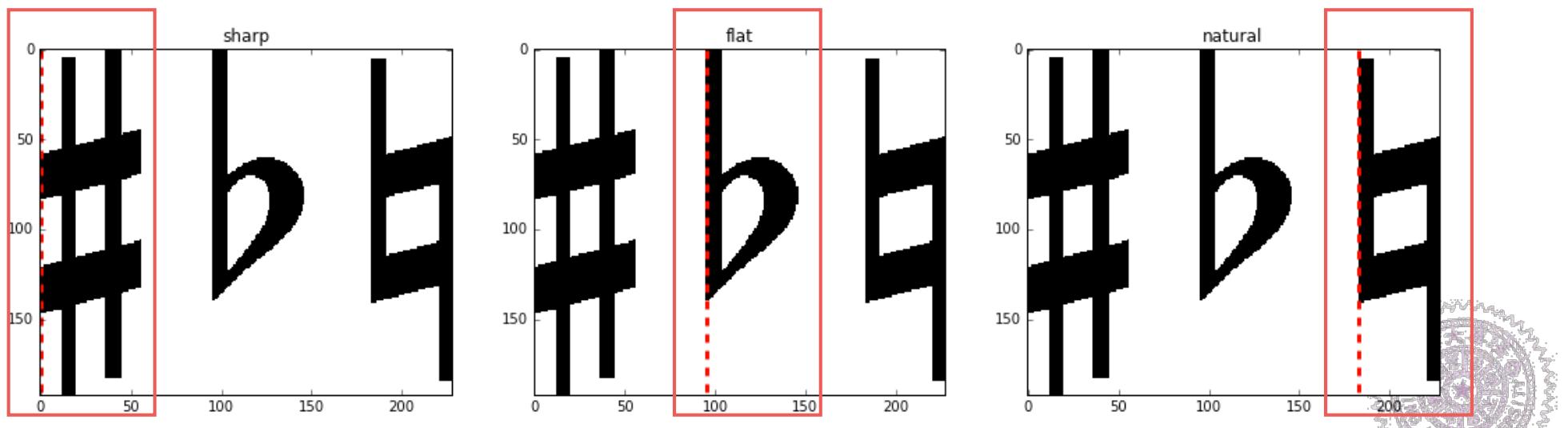
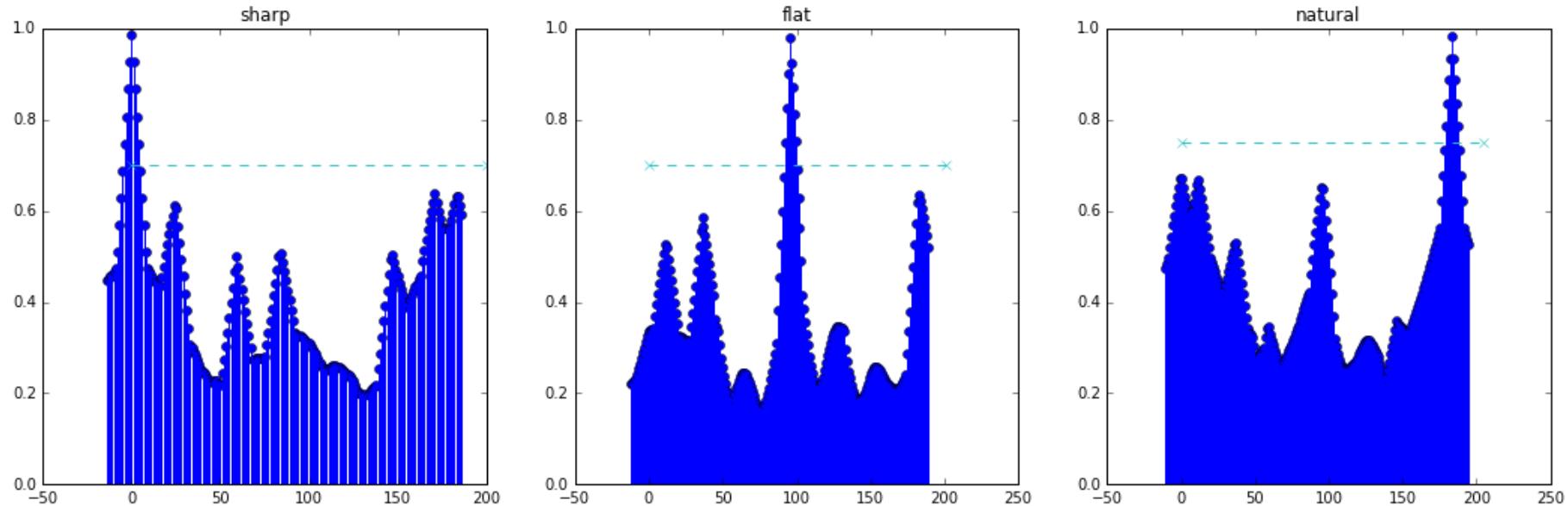


the margin between
adjacent stafflines

pitch



Accidental Detection Template Matching

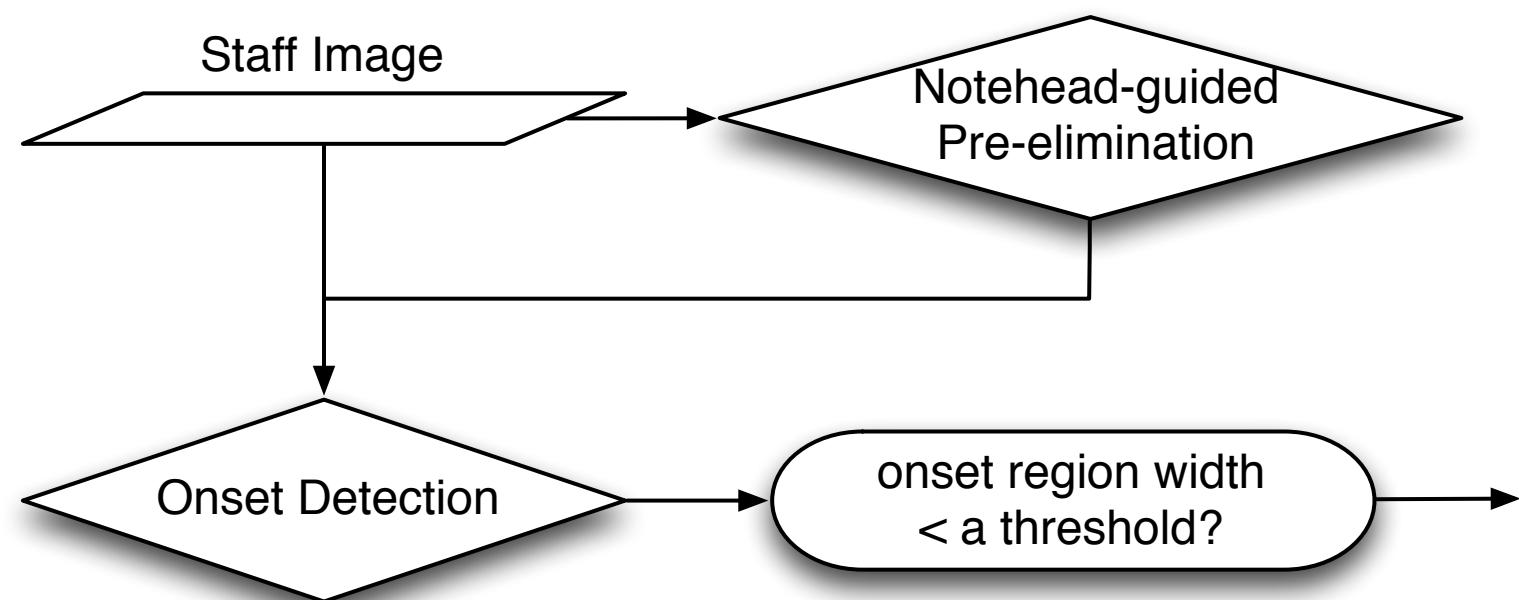
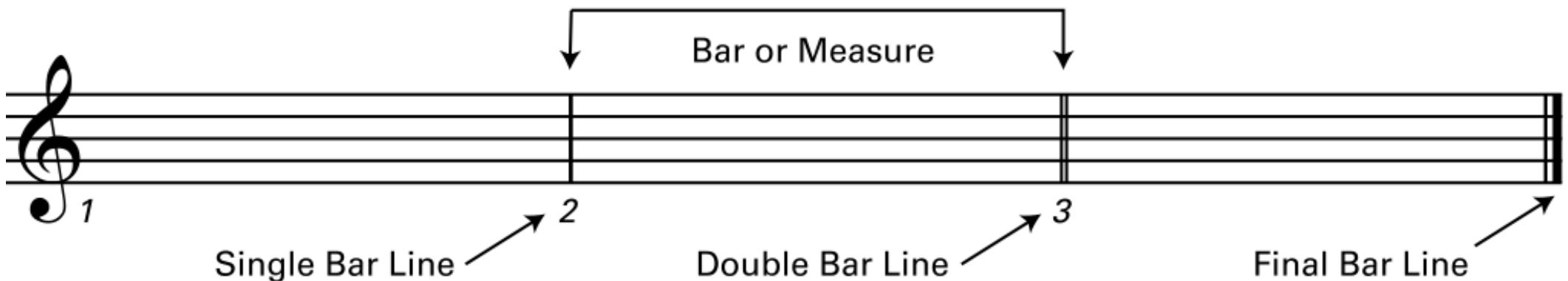


Recognition Other Symbols

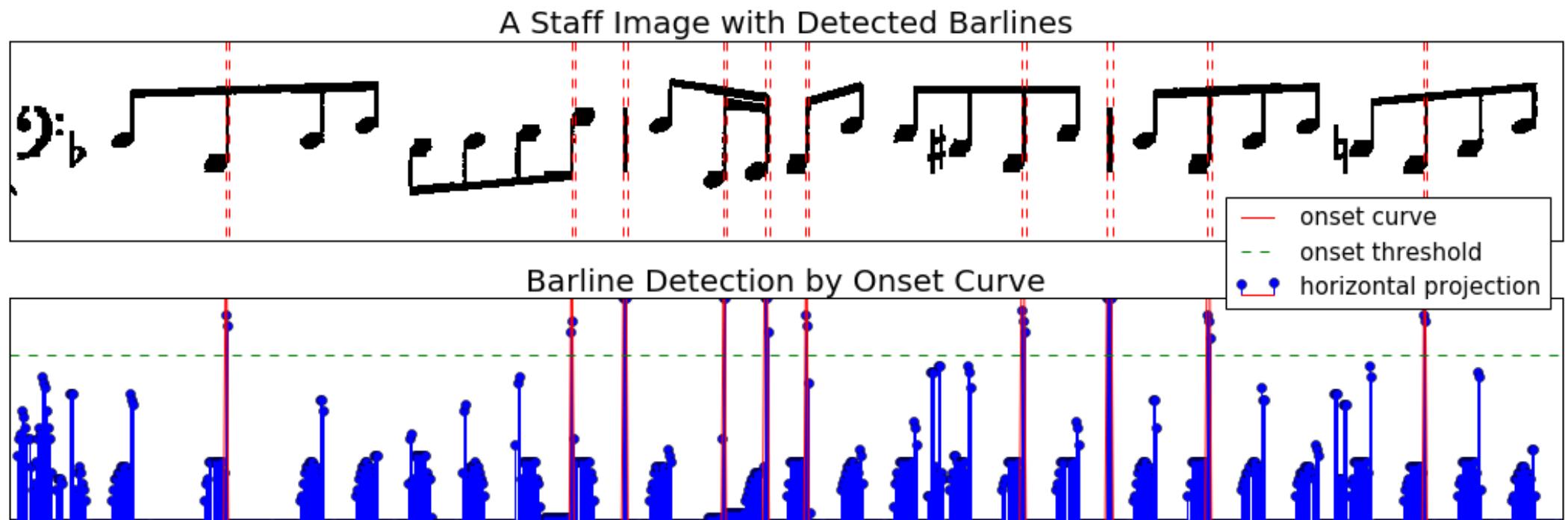


Acoustic and Hearing Group

Barline Detection

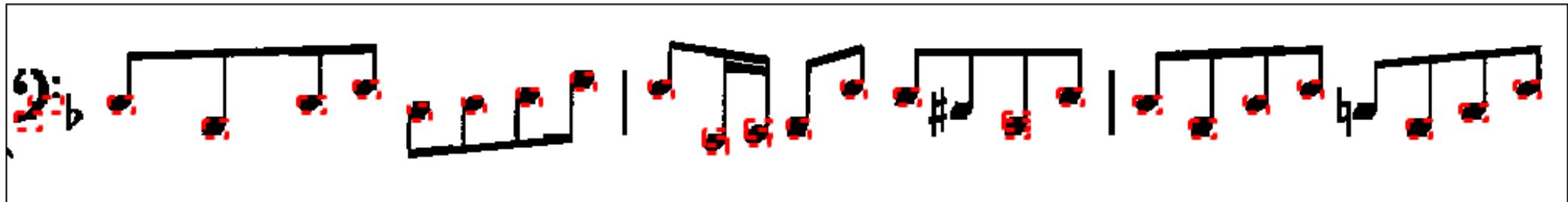


Barline Detection

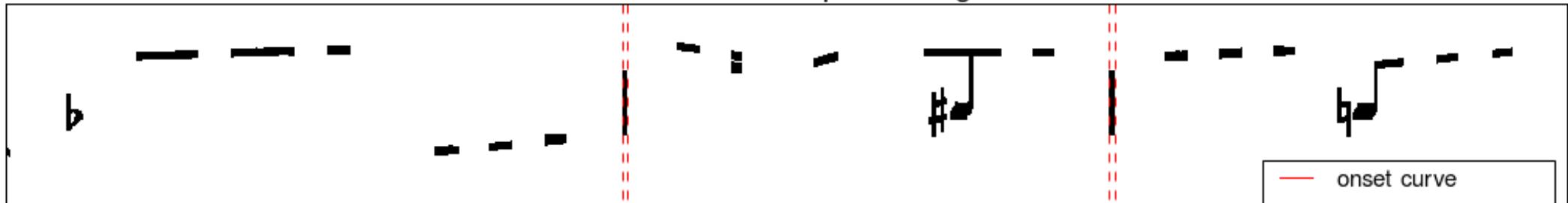


Barline Detection with Preprocessing

A Staff with Notehead Detection



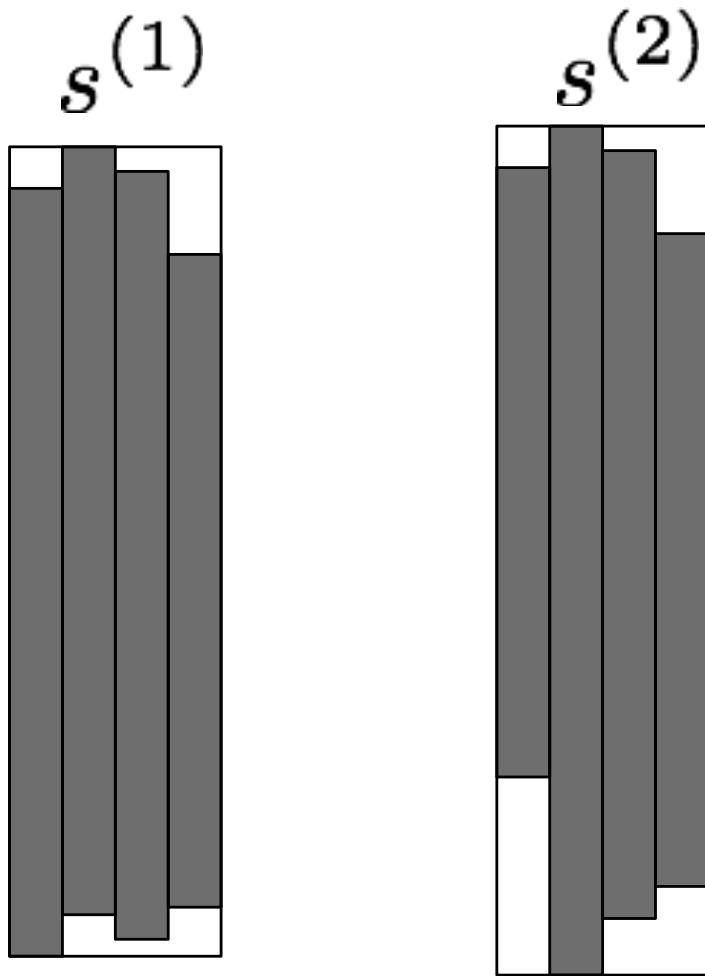
A Staff After Preprocessing



Barline Detection with Onset Curve



Stem Detection



Algorithm 1 Procedure of finding vertical lines

Input: A gray-scaled image I

Output: A set of 4-tuples

```
1:  $\mathcal{V} \leftarrow []$ 
2: for each  $c \leftarrow$  the  $x^{\text{th}}$  column of  $I$  do
3:   if  $c(y) = 1$ ,  $\forall y = y_{\min}, y_{\min} + 1, \dots, y_{\max}$  and  $y_{\max} - y_{\min} \geq 2(w + s)$  then
4:     Append  $(x, y_{\min}, y_{\max})$  to  $\mathcal{V}$ 
5:   end if
6: end for
7:  $\mathcal{R} \leftarrow []$ 
8: for each  $(x, y_{\min}, y_{\max}) \in \mathcal{V}$  do
9:   if  $\mathcal{R}$  is empty then
10:    Append  $(x, x, y_{\min}, y_{\max})$  to  $\mathcal{R}$ 
11:    Continue
12:   end if
13:    $r \leftarrow$  the reference of the last element of  $\mathcal{R}$ 
14:   if  $r_2 = x - 1$  then
15:     Update  $r_2 \leftarrow x$ 
16:     Update  $r_3 \leftarrow \min(r_3, y_{\min})$ 
17:     Update  $r_4 \leftarrow \max(r_4, y_{\max})$ 
18:   else
19:     Append  $(x, x, y_{\min}, y_{\max})$  to  $\mathcal{R}$ 
20:   end if
21: end for
22: Return  $\mathcal{R}$ 
```



Beam Grouping

Algorithm 2 Procedure of beam grouping

Input: the image I and the matrix M_{vl}

Output: A set of (\mathcal{S}, δ) 's where \mathcal{S} is the collection of stems that belong to the same

beam and $\delta = \begin{cases} 0, & \text{the beam runs through } y_{\min} \text{'s} \\ 1, & \text{the beam runs through } y_{\max} \text{'s} \end{cases}$ is a binary symbol referring to
the beam is upper or lower.

```

1:  $\mathcal{G} \leftarrow []$ 
2: for  $i \leftarrow 1, \dots, |\mathcal{R}_{vl}|$  do
3:    $s^{(i)} \leftarrow$  the  $i^{\text{th}}$  column of  $M_{vl}$ 
4:   for  $j \leftarrow i + 1, \dots, |\mathcal{R}_{vl}|$  do
5:      $s^{(j)} \leftarrow$  the  $j^{\text{th}}$  column of  $M_{vl}$ 
6:      $p_1 \leftarrow (s_1^{(i)}, s_3^{(i)}), p_2 \leftarrow (s_1^{(j)}, s_3^{(j)})$ 
7:     if  $L_I(p_1, p_2) \geq 0.9(p_{2x} - p_{1x})$  then
8:       Append  $(i, j, 0)$  to  $\mathcal{G}$ 
9:       Break
10:      end if
11:       $p_1 \leftarrow (s_1^{(i)}, s_4^{(i)}), p_2 \leftarrow (s_1^{(j)}, s_4^{(j)})$ 
12:      if  $L_I(p_1, p_2) \geq 0.9(p_{2x} - p_{1x})$  then
13:        Append  $(i, j, 1)$  to  $\mathcal{G}$ 
14:        Break
15:      end if
16:    end for
17:  end for

```

$$l_{p_1, p_2}(x) = y,$$

$$L_I(p_1, p_2) = \sum_{x=p_{1x}}^{p_{2x}} I(x, l_{p_1, p_2}(x)), p_{1x} \leq p_{2x}$$

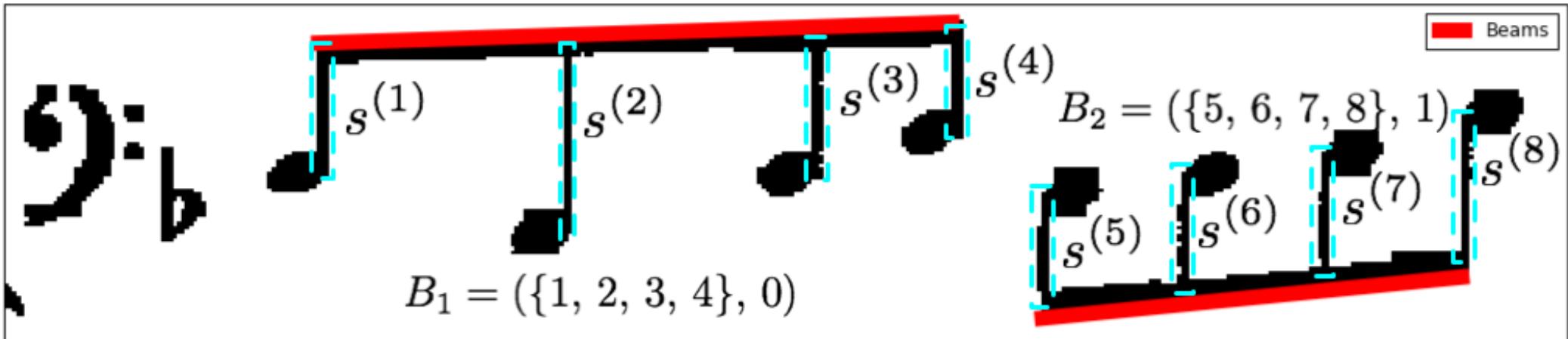
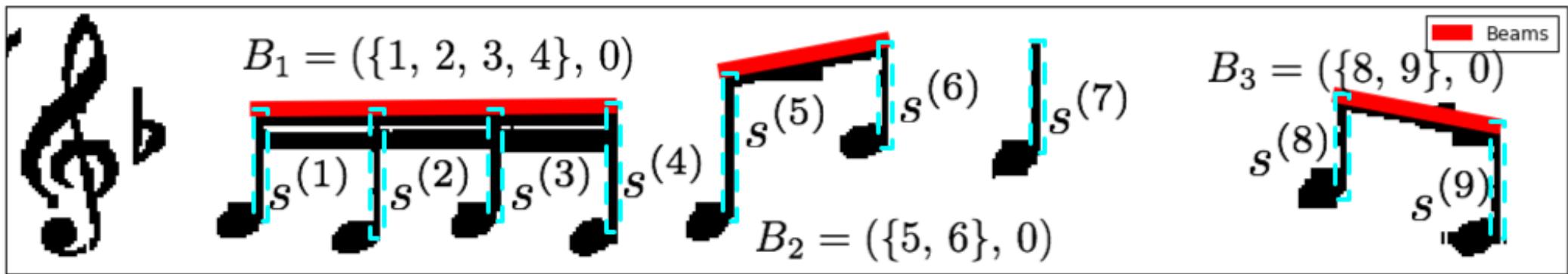
```

18:  $\mathcal{S}_{\text{beam}} \leftarrow []$ 
19: for each  $(i, j, \delta) \in \mathcal{G}$  do
20:   if  $\mathcal{S}_{\text{beam}}$  is empty then
21:     Append  $(\{i, j\}, \delta)$  to  $\mathcal{S}_{\text{beam}}$ 
22:     Continue
23:   end if
24:    $(\mathcal{S}_{\text{last}}, \delta_{\text{last}}) \leftarrow$  the last element of  $\mathcal{S}_{\text{beam}}$ 
25:    $j_{\text{last}} \leftarrow \max(\mathcal{S}_{\text{last}})$ 
26:   if  $j_{\text{last}} = i$  then
27:     Add  $j$  to  $\mathcal{S}_{\text{last}}$ 
28:     Update the last element of  $\mathcal{S}_{\text{beam}}$ 
29:   else
30:     Append  $(\{i, j\}, \delta)$  to  $\mathcal{S}_{\text{beam}}$ 
31:   end if
32: end for
33: Return  $\mathcal{S}_{\text{beam}}$ 

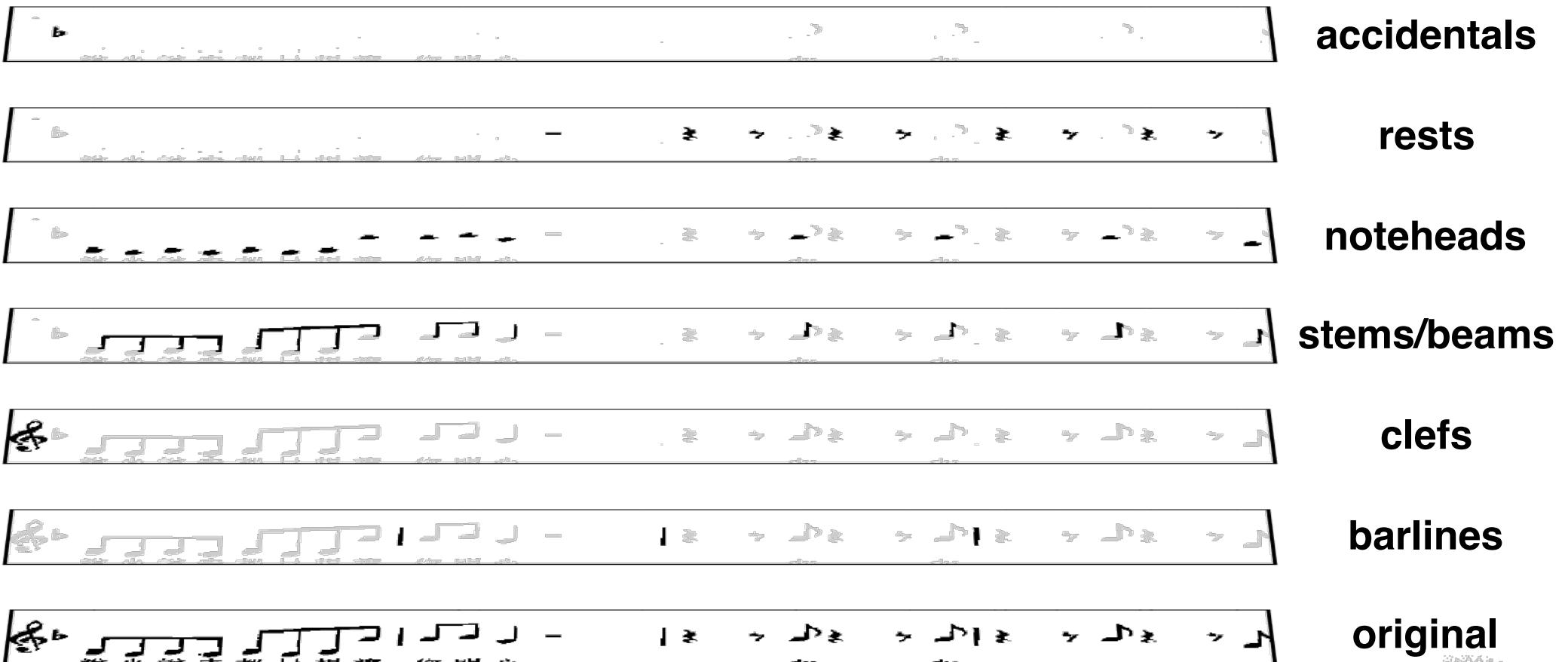
```



Beam Grouping



Hierarchical Detection



Summary of Detection on Different Symbols

	Detection Approach	Locating Symbols	Parameters	Working Images
Barlines	Horizontal Projection	Onset Curve	width, onset threshold	Staff
Clefs	Double-stage SVM	Onset Curve	merge, onset threshold	Measure
Rests	Template Matching	Onset Curve	similarity, onset threshold	Measure
Accidentals	Template Matching	Sliding Window	similarity threshold	Pitch
Noteheads	Double-stage SVM	Onset Curve	onset threshold	Pitch
		Stem-guided	N/A	Measure
Stems	RLE	Sliding Window	length threshold	Measure
Beams	Line Integral	Stem-guided	line integral threshold	Measure



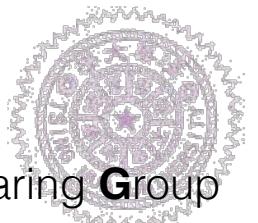
Summary of Features on Different Symbols

	Features
Barlines	RLE of Horizontal Projection
Clefs	Tiny Image, Horizontal/Vertical Projection and Centroid
Rests	Cosine Similarity
Accidentals	Cosine Similarity
Noteheads	Tiny Image, Diagonal Projection
Stems	RLE
Beams	Rectangles of Stems



References

- [1] T. Pinto, A. Rebelo, G. Giraldi, and J. S. Cardoso, “Music score binarization based on domain knowledge,” *Pattern Recognition and Image Analysis - 5th Iberian Conf. (IbPRIA)*, pp. 700–708, 2011.
- [2] O. Nobuyuki, “A threshold selection method from gray-level histograms,” *IEEE Trans. Systems, Man and Cybernetics*, vol. 9, pp. 62–66, 1979.
- [3] Q. Chen, Q.-s. Sun, P. A. Heng, and D.-s. Xia, “A double-threshold image binarization method based on edge detector,” *Pattern Recognition*, vol. 41, pp. 1254–1267, 2008.
- [4] L.-K. Huang and M.-J. J. Wang, “Image thresholding by minimizing the measures of fuzziness,” *Pattern Recognition*, vol. 28, pp. 41–51, 1995.
- [5] D.-M. Tsai, “A fast thresholding selection procedure for multimodal and unimodal histograms,” *Pattern Recognition Letters*, vol. 16, pp. 653–666, 1995.
- [6] J. Bernsen, “Dynamic thresholding of grey-level images,” in *Proc. the 8th. Int. IEEE Conf. CAD Systems in Microelectronics (CADSM)*, pp. 1254–1267, 2005.
- [7] R. Randriamahefa, J. P. Cocquerez, F. Fluhr, C. Pepin, and S. Philipp, “Printed music recognition,” in *Proc. the Second Int. Conf. on Document Analysis and Recognition*, pp. 898–901, 1993.
- [8] K. T. Reed and J. Parker, “Automatic computer recognition of printed music,” in *Proc. the 13th Int. Conf. on Pattern Recognition*, vol. 3, p. 803–807, 1996.
- [9] P. Bellini, I. Bruno, and P. Nesi, “Optical music sheet segmentation,” in *Proc. the First Int. Conf. on Web Delivering of Music*, pp. 183–190, 2001.
- [10] H. Miyao, “Stave extraction for printed music scores,” *Intelligent Data Engineering and Automated Learning-IDEAL 2002*, H. Yin, N. Allinson, R. Freeman, J. Keane, and S. Hubbard, Eds. Springer, pp. 621–634, 2002.
- [11] F. Rossant and I. Bloch, “Robust and adaptive OMR system including fuzzy modeling, fusion of musical rules, and possible error detection,” *EURASIP J. on Advances in Signal Processing*, vol. 1, 2007.
- [12] C. Dalitz, M. Droettboom, B. Pranzas, and I. Fujinaga, “A comparative study of staff removal algorithms,” *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 30, pp. 753–766, 2008.



References

- [13] J. d. S. Cardoso, A. Capela, A. Rebelo, C. Guedes, and J. Pinto da Costa, "Staff detection with stable paths," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 31, p. 1134–1139, 2009.
- [14] A. Dutta, U. Pal, A. Fornés, and J. Lladós, "An efficient staff removal approach from printed musical documents," in *Proc. the 20th Int. Conf. on Pattern Recognition*, p. 1965–1968, IEEE Computer Society, 2010.
- [15] A. Rebelo, G. Capela, and J. S. Cardoso, "Optical recognition of music symbols: a comparative study," *Int. J. Document Analysis Recognition*, vol. 13, pp. 19–31, 2010.
- [16] I. Leplumey, J. Camillerapp, and G. Lorette, "A robust detector for music staves," in *Proc. the Int. Conf. on Document Analysis and Recognition*, p. 203–210, 1993.
- [17] Q.-A. Arshad, W. Z. Khan, and Z. Ihsan, "Overview of algorithms and techniques for optical music recognition."
- [18] I. Fujinaga, *Adaptive Optical Music Recognition*. PhD thesis, Faculty of Music, McGill University, Montréal, Canada, 1996.
- [19] A. Forns, A. Dutta, A. Gordo, and J. Llads, "CVC-MUSCIMA: A ground-truth of handwritten music score images for writer identification and staff removal," *Int. J. on Document Analysis and Recognition*, vol. 15, pp. 243–251, 2012.
- [20] T. Kanungo, R. M. Haralick, H. S. Baird, W. Stuezle, and D. Madigan, "A statistical, nonparametric methodology for document degradation model validation," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 22, pp. 1209–1223, 2000.
- [21] L. Pugin, "Optical music recognition of early typographic prints using hidden Markov models," in *Proc. the Int. Society for Music Information Retrieval*, pp. 53– 56, 2006.
- [22] A. Fornés, S. Escalera, J. L Ladòs, G. Sàncchez, P. Radeva, and O. Pujol, "Handwritten symbol recognition by a boosted blurred shape model with error correction," in *Proc. the 3rd Iberian Conf. on Pattern Recognition and Image Analysis, Part I*. Springer, Berlin, pp. 13–21, 2007.
- [23] H. Miyao and Y. Nakano, "Head and stem extraction from printed music scores using a neural network approach," in *Proc. Third Int. Conf. on Document Analysis and Recognition*, pp. 1074–1079, IEEE Computer Society, 1995.
- [24] L. Chen and C. Raphael, "Human-directed optical music recognition," in *Electronic Imaging, Document Recognition and Retrieval XXIII*, pp. 1–9, Society for Imaging Science and Technology, 2016.
- [25] M. A. Fischler and R. C. Bolles, "Random sample consensus: A paradigm for model fitting applications to image analysis and automated cartography," in *Proc. Image Understanding Workshop*, pp. 71–88, 1980.
- [26] C.-W. Hsu, C.-C. Chang, and C.-J. Lin, *A Practical Guide to Support Vector Classification*. National Taiwan University, Taipei 106, Taiwan, 2016.

