
Training Neural Nets to Achieve Audio-to-Score Translation: Opening the Black-Box

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 It is suggested that the task of audio-to-score translation offers an adequate testbed
2 to investigate the division of labor between background knowledge and machine
3 learning in the domain of audio pattern recognition, with a controllable level of
4 difficulty and the ability to synthesize a limitless amount of labelled data.

5 As a proof of concept, this paper focuses on pitch detection from audio signals.
6 Extensive background knowledge is used to initialize simple convolutional neural
7 nets (NN) and achieve the recognition of single notes with a decent accuracy. The
8 performance achieved by trained NNs, however, is significantly higher. Some
9 tentative interpretations of this fact are obtained by opening the black box and
10 inspecting the modifications of the NN filters due to supervised learning.

11 1 Introduction and Rationale

12 Audio-to-speech translation has been intensively studied since the early 80s, as a key step of the
13 natural language processing chain. While mainstream approaches used to be rooted in statistical
14 and linguistic feature extraction [1], new end-to-end approaches rooted on deep neural networks are
15 revolutionizing the field [2], along the same lines as in computer vision [3, 4].

16 This paper investigates a related but different issue, the audio-to-score translation (AST) in music, at
17 the core of the Music Information Retrieval Evaluation eXchange (MIREX) international challenges
18 [5]. This domain offers rich resources in terms of labelled data and feature extraction using extensive
19 signal processing libraries. It thus supports the investigation of the benefits of using background
20 knowledge, and specifically what knowledge can be given to the learning agent and what should
21 better be learned, in a systematic way.

22 **Data resources.** MIREX 2018 offers a number of tracks, ranging from the identification of the
23 type of music, to the name of the composer or the music mood. Former MIREX editions were
24 concerned with other tracks including the alignment of audio to scores [6]. In the following, we
25 focus on a subtask of audio-to-score translation task, *pitch detection*, formulated as a supervised
26 learning task and tackled using the MAPS database [7]. It is claimed that the single note identification
27 from audio signal is most similar to phoneme identification from audio signal, the elementary task of
28 audio-to-speech translation.

29 **Prior knowledge and pre-processing** Dedicated and robust signal processing (SP) libraries[8]
30 have been developed to achieve audio signal pre-processing and build time-frequency representations
31 of the signal. These signal representations can be provided as pixel matrices or vectors to standard
32 neural network (NN) architectures, together with the associated labels (single note or chords), and
33 used to train NNs using standard supervised learning.

34 An alternative is offered by using simple convolutional NN architectures on the raw signal. The
35 point is that the convolutional mask aimed to detect a given frequency f can be either learned, or
36 deterministically set to $W_{f,j} = \cos(\frac{2\pi f j}{f_s})$ for $j = [[0, N]]$, with N is the mask dimension and
37 $f_s = 22.05$ kHz the sampling frequency of the audio signal.

38 **Research Agenda** This paper reports on the lessons learned from training NN (either from scratch,
39 or using educated weight initialization) for pitch detection, claiming that these lessons are relevant to
40 audio-to-speech phoneme recognition for the following reasons.

41 Pitch detection is assumedly easier than phoneme identification for a ground truth compact description
42 of every note is available, e.g. in terms of Fourier coefficient. The diversity among notes is lesser than
43 for phonemes (a note most similar but different from another note does not exist as far as musical
44 instruments are well tuned), but still presents some ambiguities as a note evokes its harmonic by
45 construction. The phoneme diversity due to the different speakers has its equivalent in music: a same
46 note played on different instruments corresponds to different audio signals, due to the distinctive
47 timbre of the musical instruments.

48 For these reasons, it is conjectured that pitch detection defines a relevant and informative testbed for
49 phoneme detection. After briefly presenting Automatic Music Transcription, we describe the goals
50 of the presented experiments, before reporting on the lessons learned, about how and when the SP
51 knowledge is worth being exploited and what are the pitfalls.

52 2 Automatic Music Transcription

53 Automatic Music Transcription (AMT) [9] is a difficult task even for human experts [10]. Most of
54 the time, AMT and pitch detection research is done on piano music, to avoid the variability due
55 to the direct contact between the musician and the string that produces the sound. First attempts
56 to automatically transcribe polyphonic music due to Moorer [11] inspired many followers [12, 13,
57 14], mostly using statistical approaches. AMT and Music Information Retrieval were more recently
58 revisited using Deep Learning techniques, specifically convolutional (CNNs) and deep belief networks
59 (DBNs) [15, 16, 17].

60 Only pitch detection is considered in the following. The set of audio frequencies corresponding to
61 the musical notes are distributed on a log-scale; the pythagorean scale consists in $f_n \propto f_0 2^{n/12}$,
62 mapping the human ear perception.

63 Pitch detection can be achieved through extracting robust features from the raw audio signal. Com-
64 pared to computer vision, the extracted features are only required to be invariant w.r.t. translation in
65 time.¹ In principle, each such feature would correspond to a given pitch, facilitating the interpretation
66 of the feature extraction achieved by the neural net.

67 3 Experimental setting and goals of experiments

68 The main goal of these experiments regards the nature and value of the prior SP knowledge concerning
69 the target identification task. Specifically, does the available prior knowledge yield a perfect accuracy?
70 Otherwise, how is this knowledge modified and improved through supervised training? Thirdly, how
71 much gain is provided by the use of the prior knowledge? Finally, how do the NN hyper-parameters
72 interact with the model space and how to best tune them?

73 These questions are empirically investigated on the MIDI Aligned Piano Sounds dataset [7]. We
74 restricted ourselves to the identification of monophonic (isolated) notes, forming 88 distinct balanced
75 classes. The raw audio signal is downsampled to 22kHz. Each file is divided into 2,400 sample-long
76 examples (about 109ms), overlapping with a stride of 600. The period of the lowest note (A0,
77 frequency 27.5Hz) spreads over 800 samples.

78 The NN architecture is a convolutional neural network applied on the raw audio signal, with a stride
79 of 10 or 50, followed by a global max pooling. Each note is associated to one filter mask of dimension

¹In all generality, the extracted features should also be invariant w.r.t. the timbre of the instruments; however only piano music is considered in the following.

80 800. The educated initialization of the corresponding filter weights is defined from the pure sine wave
 81 corresponding to the fundamental frequency (with no biases for the sake of interpretability).

82 Three architectures are compared: **model A** is a 1 layer CNN model with *initialized weights*, with
 83 *no training*, **model B** a 1 layer CNN model with *initialized weights*, with *training*, and **model C** a
 84 1-layer CNN model with *random initialization*, with *training*.

85 Model performance is assessed using 8-fold cross validation (all samples of a same audio file being
 86 either in the training or in the test set).

87 4 Results

88 **Impact of prior knowledge on accuracy and learning curve.** As shown (Table 1), untrained
 89 model A yields 45% accuracy (significantly higher than the default accuracy for 88 balanced classes),
 90 thanks to its educated initialization. Still, trained models B and C both significantly outperform model
 91 A. The educated initialization (in model B) speeds up the convergence (compared to model C) toward
 92 the same optimum (Fig. 1). The performances are improved as the stride is reduced, especially so for
 93 high frequencies (Fig. 3).

Table 1: Comparison of the performances of models A, B and C on test sets after 100 epochs.

| | Stride 50 | | Stride 10 | |
|-----------------|-----------------|-----------------|-----------------|-----------------|
| | Before training | After training | Before training | After training |
| initialized | 0.45 | 0.64 ± 0.02 | 0.46 | 0.71 ± 0.02 |
| not initialized | - | 0.64 ± 0.02 | - | 0.70 ± 0.01 |

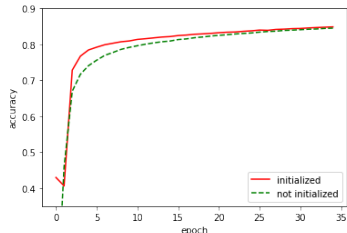


Figure 1: Learning curves of trained Model B (initialized with SP knowledge) and Model C (random initialization) on the training set (stride 50).

94 **Opening the black-box.** The masks in models A, B and C are inspected by plotting their Fourier
 95 coefficients (Fig. 2). By construction, masks in model A (Fig. 2.top) exactly correspond to the
 96 sought frequencies. On model B (SP initialization), supervised learning results in augmenting
 97 each filter with the harmonics of the associated frequency, particularly so for low frequencies; this
 98 augmentation explains the improved performance of model B compared to model A. On model
 99 C (random initialization), the same augmentation is observed with a notable difference: for low
 100 frequencies f , only harmonics $2f, 3f, \dots$ are learned; frequency f itself is missed. This phenomenon
 101 confirms that good performances *per se* do not imply that the targeted concepts have been correctly
 102 identified.

103 The accuracy per note of models A, B, and C are depicted on Fig. 3. As mentioned above, we observe
 104 a significant improvement in performances when the model is trained (B and C), above all in a low
 105 frequency notes range. Indeed, for pitches between 40 and 80, the accuracy is roughly equivalent,
 106 whereas for pitches < 40 , the trained masks obtain much better results. This is probably due to the
 107 timbral specificity of the piano, where the distribution of the energy for a given note is distributed
 108 differently depending on the pitch range.

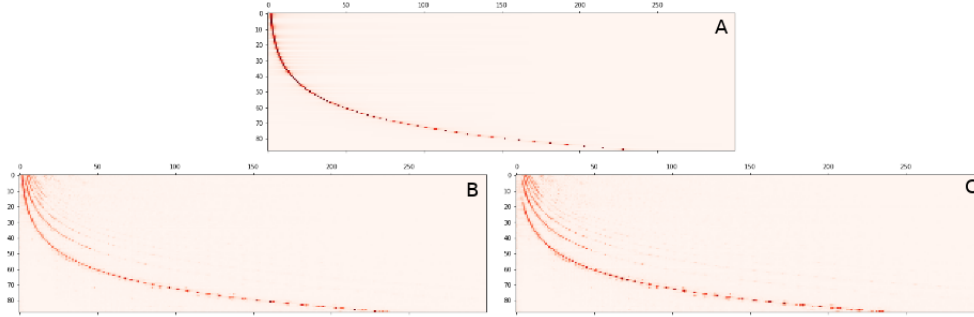


Figure 2: Models A, B and C: magnitude of frequencies detected for each filter (horizontal line), with stride 50, after 100 epochs training for models B and C.

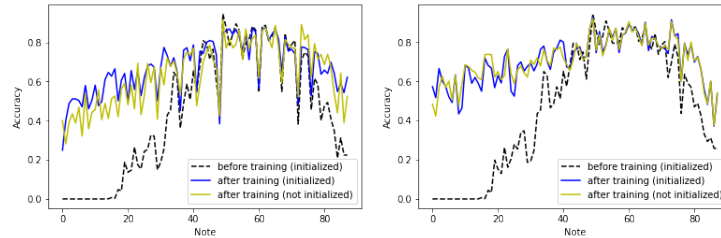


Figure 3: Models A, B and C: Accuracy per note. Left: stride 50; Right: stride 10. The downward accuracy spikes occur for frequencies whose period is a divider of the stride period (see discussion).

109 **Adding a fully connected layer.** Another model is investigated, adding a fully connected layer on
 110 the top of the conv/max pooling layers. The question is whether and how this additional FC layer
 111 will alleviate the limitations of the simple B and C architectures, and in practice, whether augmenting
 112 the filters with the associated harmonics will still be necessary. Most interestingly, it appears that the
 113 augmentation of the filters is no longer necessary (Fig. 4, left) and the fact that both a frequency and
 114 its harmonics are involved in the pitch detection is visible from the FC weights, below the diagonal
 115 (Fig. 4, right). This phenomenon is observed mostly with educated initialization. For both educated
 116 and random initialization, the predictive accuracy increases from circa 70% to 80%.

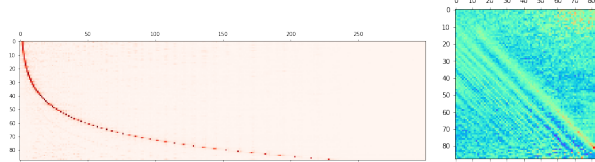


Figure 4: Model B augmented with a fully connected layer. Left: magnitude of frequencies detected for each filter; Right: weights of the fully connected layer; pixel (i,j) represents the weight of filter f_j to identify frequency f_i .

117 5 Discussion

118 These experiments firstly confirm that good performances can be obtained through merely extracting
 119 features correlated to the target class (here the harmonics of the sought frequencies). The convolutional
 120 structure, aimed to achieve phase-invariance, is most effective for low stride values. For higher stride
 121 values, the invariance property does not hold and the accuracy shows a downward peak. An alternative
 122 is to consider neurons in the complex space [18].

123 Overall, the merits of using SP knowledge are twofold: it speeds up the convergence and it ensures
 124 that the target concept is properly grasped. On the other hand, supervised learning is beneficial as it
 125 automatically detects and exploits the relationships among frequencies and harmonics for a more
 126 robust detection.

127 **References**

- 128 [1] S. J. Godsill et al. “Bayesian computational methods for sparse audio and music processing”.
 129 In: *2007 15th European Signal Processing Conference*. Sept. 2007, pp. 345–349.
- 130 [2] Aäron van den Oord et al. “WaveNet: A Generative Model for Raw Audio”. In: *CoRR*
 131 abs/1609.03499 (2016). arXiv: 1609.03499. URL: <http://arxiv.org/abs/1609.03499>.
- 132 [3] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. The MIT Press, 2016.
 133 ISBN: 0262035618, 9780262035613.
- 134 [4] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. “ImageNet Classification with Deep
 135 Convolutional Neural Networks”. In: *Proceedings of the 25th International Conference on*
 136 *Neural Information Processing Systems - Volume 1. NIPS’12*. Lake Tahoe, Nevada: Curran
 137 Associates Inc., 2012, pp. 1097–1105. URL: [http://dl.acm.org/citation.cfm?id=](http://dl.acm.org/citation.cfm?id=2999134.2999257)
 138 [2999134.2999257](http://dl.acm.org/citation.cfm?id=2999134.2999257).
- 139 [5] *Mirex Challenge*. URL: https://www.music-ir.org/mirex/wiki/MIREX_HOME.
- 140 [6] Damien Garreau et al. “Metric Learning for Temporal Sequence Alignment”. In: *Advances in*
 141 *Neural Information Processing Systems 27*. Ed. by Z. Ghahramani et al. Curran Associates, Inc.,
 142 2014, pp. 1817–1825. URL: [http://papers.nips.cc/paper/5383-metric-learning-](http://papers.nips.cc/paper/5383-metric-learning-for-temporal-sequence-alignment.pdf)
 143 [for-temporal-sequence-alignment.pdf](http://papers.nips.cc/paper/5383-metric-learning-for-temporal-sequence-alignment.pdf).
- 144 [7] Valentin Emiya, Roland Badeau, and Bertrand David. “Multipitch Estimation of Piano Sounds
 145 Using a New Probabilistic Spectral Smoothness Principle”. In: *IEEE Trans. Audio, Speech*
 146 *& Language Processing* 18.6 (2010), pp. 1643–1654. DOI: 10.1109/TASL.2009.2038819.
 147 URL: <https://doi.org/10.1109/TASL.2009.2038819>.
- 148 [8] *Librosa library*. URL: <https://librosa.github.io/librosa/>.
- 149 [9] Emmanouil Benetos et al. “Automatic music transcription: challenges and future directions”.
 150 In: *J. Intell. Inf. Syst.* 41.3 (2013), pp. 407–434. DOI: 10.1007/s10844-013-0258-3. URL:
 151 <https://doi.org/10.1007/s10844-013-0258-3>.
- 152 [10] George Tzanetakis. “Anssi Klapuri, Manuel Davy, Eds: Signal Processing Methods for Music
 153 Transcription”. In: *Computer Music Journal* 32.4 (2008), pp. 86–88. DOI: 10.1162/comj.
 154 2008.32.4.86. URL: [https://doi.org/10.1162/comj.](https://doi.org/10.1162/comj.2008.32.4.86)
 155 [2008.32.4.86](https://doi.org/10.1162/comj.2008.32.4.86).
- 155 [11] James A. Moorer. “On the Transcription of Musical Sound by Computer”. In: *Computer Music*
 156 *Journal* 1.4 (1977), pp. 32–38. ISSN: 01489267, 15315169. URL: [http://www.jstor.org/](http://www.jstor.org/stable/40731298)
 157 [stable/40731298](http://www.jstor.org/stable/40731298).
- 158 [12] Keith D. Martin. *A Blackboard System for Automatic Transcription of Simple Polyphonic*
 159 *Music*. Tech. rep. 1996.
- 160 [13] M. Privosnik and M. Marolt. “A system for automatic transcription of music based on multiple-
 161 agents architecture”. In: *MELECON ’98. 9th Mediterranean Electrotechnical Conference.*
 162 *Proceedings (Cat. No.98CH36056)*. Vol. 1. May 1998, 169–172 vol.1. DOI: 10.1109/MELCON.
 163 1998.692363.
- 164 [14] R. Keren, Y. Y. Zeevi, and D. Chazan. “Multiresolution time-frequency analysis of polyphonic
 165 music”. In: *Proceedings of the IEEE-SP International Symposium on Time-Frequency and*
 166 *Time-Scale Analysis (Cat. No.98TH8380)*. Oct. 1998, pp. 565–568. DOI: 10.1109/TFSA.
 167 1998.721487.
- 168 [15] Eric J. Humphrey, Juan P. Bello, and Yann Lecun. “Feature Learning and Deep Architectures:
 169 New Directions for Music Informatics”. In: *J. Intell. Inf. Syst.* 41.3 (Dec. 2013), pp. 461–481.
 170 ISSN: 0925-9902. DOI: 10.1007/s10844-013-0248-5. URL: [http://dx.doi.org/10.](http://dx.doi.org/10.1007/s10844-013-0248-5)
 171 [1007/s10844-013-0248-5](http://dx.doi.org/10.1007/s10844-013-0248-5).
- 172 [16] E. J. Humphrey, T. Cho, and J. P. Bello. “Learning a robust Tonnetz-space transform for
 173 automatic chord recognition”. In: *2012 IEEE International Conference on Acoustics, Speech*
 174 *and Signal Processing (ICASSP)*. Mar. 2012, pp. 453–456. DOI: 10.1109/ICASSP.2012.
 175 6287914.
- 176 [17] Siddharth Sigtia, Emmanouil Benetos, and Simon Dixon. “An End-to-End Neural Network for
 177 Polyphonic Piano Music Transcription”. In: *IEEE/ACM Trans. Audio, Speech & Language*
 178 *Processing* 24.5 (2016), pp. 927–939. DOI: 10.1109/TASLP.2016.2533858. URL: [https://doi.org/10.](https://doi.org/10.1109/TASLP.2016.2533858)
 179 [1109/TASLP.2016.2533858](https://doi.org/10.1109/TASLP.2016.2533858).
- 180 [18] Chiheb Trabelsi et al. “Deep Complex Networks”. In: *CoRR* abs/1705.09792 (2017). arXiv:
 181 1705.09792. URL: <http://arxiv.org/abs/1705.09792>.