

PHONEMATIC TRANSLATION OF POLISH TEXTS BY THE NEURAL NETWORK

A. BIELECKI, I.T. PODOLAK, J. WOSIEK*

Institute of Computer Science, Jagellonian University
Nawojki 11, 30-072 Kraków, Poland

AND

E. MAJKUT

Institute of Polish Philology, Jagellonian University
Gołębia 14-20, 31-007 Kraków, Poland

(Received March 26, 1996)

Using the backpropagation algorithm, we have trained the feed forward neural network to pronounce Polish language, more precisely to translate Polish text into its phonematic counterpart. Depending on the input coding and network architecture, 88%-95% translation efficiency was achieved.

PACS numbers: 89.80.+h

1. Introduction

Computer science has been interested in phonematic translations for a long time. So far the problem was approached in two different ways. A simple implementation of all transformation rules, including irregularities, is the first possibility. A DECTalk system converting English text to speech was implemented in such a way [1]. In such commercial systems a look-up table of about million bits is used to store the phonetic transcription of common and irregular words, and phonological rules are applied to words that are not in this table [2, 3]. A similar system for Polish is described in [4]. Using an artificial neural network is the second possibility. A NETTalk, described in [5] and [6], and also mentioned in [7] where its abilities and

* Supported in part by the KBN grants PB 2P03B19609 and PB 2P30225206.

implementation difficulties in comparison with the DECtalk are discussed, is a network system that learns to convert English text to speech.

In this paper we describe few neural networks that learn to pronounce Polish text. Each of them has slightly different architecture and differs with respect to others, both in learning speed and translation efficiency.

2. Phonemes and phonematic transformations in Polish

Polish language is much simpler than English to specify all transformation rules. Bolc and Maksymienko have implemented a system based on this approach and have obtained good results [4]. We have used an artificial neural network for this purpose in order to test both approaches and also to compare results with those of Sejnowski's NETtalk implementation for the English language.

The implemented neural networks transform letter with its context (*i.e.* letters before and behind the letter) into a phoneme. A phoneme is the smallest part of a language which differentiates a meaning but has no meaning itself. For instance, two Polish words:

tam [**ta**m] (there) and

sam [**sa**m] (alone)

differ only by the parts [t] and [s]. Therefore [t] and [s] are phonemes according to the definition given above. We have defined phoneme by its functional role it plays in a word but another definition can be put forward: a phoneme is a set of features which distinguish the given phoneme from another one. Such features are called distinctive. In this way the phoneme is defined by its structure, *i.e.* the way it is pronounced. As an example consider two Polish words:

dam [**da**m] (I will give) and

tam [**ta**m] (there).

The phoneme [d] is voiced, dental, buccal, hard whereas [t] is unvoiced, dental, buccal, hard. It has been shown that a voice is a distinctive feature. We use the Jakobson theory [8] of phoneme description. In this theory all phonologic features are binary hence each phoneme is represented by a vector with binary entries. Fourteen distinctive features (binary ingredients) are needed to describe all Polish phonemes unambiguously (see Tables I to III). However not all of them are required for every phoneme.

TABLE I

The vocal phonemes with their features. Other possible features (*ie. vocal, consonantal, nasal, palatal, labial, lateral, fricative, dental, explosive, and voiced*) are not distinctive among them.

feature	u	o	i	e	a	y
velar	+	+	-	-	-	-
high-pitched	+	-	+	-	-	+
low	-	-	-	-	+	-
centralized	-	-	-	-	-	+

TABLE II

The consonantal phonemes with their features. Other possible features (*ie. vocal, consonantal, high-pitched, low, centralized, fricative, explosive, and voiced*) are not distinctive among them.

feature	ń	m	n	u	u̯	ĩ	l	r
velar	-	-	-	+	+	-	-	-
nasal	+	+	+	-	+	-	-	-
palatal	+	-	-	-	-	+	-	-
labial	-	+	-	-	-	-	-	-
lateral	-	-	-	-	-	-	+	-
dental	-	-	+	-	-	-	-	-

TABLE III

The consonantal phonemes with their features. Other possible features (*ie. vocal, consonantal, high-pitched, low, centralized, nasal, and lateral*) are not distinctive among them.

feature	v	f	b	p	x	g	k	ż	ś	g̃	k̃	ć	z	s	ż	ś	d	t	ź	c	ć	č
velar	-	-	-	-	+	+	+	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
palatal	-	-	-	-	-	-	-	+	+	+	+	+	+	-	-	-	-	-	-	-	-	-
labial	+	+	+	+	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
fricative	+	+	-	-	+	-	-	+	+	-	-	-	+	+	+	+	-	-	-	-	-	-
dental	-	-	-	-	-	-	-	-	-	-	-	-	+	+	-	-	+	+	+	+	-	-
explosive	-	-	+	+	-	+	+	-	-	+	+	-	-	-	-	-	+	+	-	-	-	-
voiced	+	-	+	-	-	+	-	+	-	+	-	+	-	+	-	+	-	+	-	+	-	+

Care must be taken while implementing automatic transformation between the text and phonematic transcription. The first problem is that the transformation of a single letter into a phoneme is usually ambiguous. In most cases it depends on the context in which the letter appears, *i.e.* it depends on the letters (or, generally, characters) both before and after the letter in question. For example, the letter sequence *dzi* gives in effect a single

phoneme [ć] and using the rules of Polish language letter *d* is transformed into [ć] and both *z* and *i* into empty phonemes [4]. Thus the transformation of *d* in this sequence depends upon two letters following it, transformation of *z* upon both one letter before and one after it, and that of *i* upon two letters preceding it. Some letters are independent of the context (like *a*, *j*, *l*, *ł*, *ń*, *o*, and *ó*), some are dependent only on the neighbouring letters, and some, yet less numerous, depend on two or three letters. In addition, an interwords and midwords assimilations exist in Polish, for instance:

kosz truskawek [koś-truskavek] (a basket full of strawberries) and
kosz gruszek [koż-grušek] (a basket full of pears). Variations of transformation of Polish nasals (ą, ę) provide another example of similar difficulties:
ręka [re ɥka] (hand): ę → e ɥ

dębu [dembu] (oak – a genetive form): ę → em

Kęty [kenty] (a name of the town in Poland): ę → en

wzięła [vźeɥa] (she gave): ę → e ɥ

kręcił [kreńciw] (he turned round): ę → eń

Additional difficulties appear in rare and unorthodox events (for instance: *zamarzać* [zamarzać] – to freeze). Transformations for which only few examples exist cause also some problems (for instance: *ź* → empty phoneme, *d* → [C] in Polish). Furthermore, there are no rules for a few borrowings and specialistic terms [9]. Moreover, the text can be pronounced correctly in different ways. There exist not only classical and common Polish but Warsaw and Cracow pronunciation as well [10]. We taught the implemented neural network according to the classical Cracow pronunciation basing on [11]. A list of various transformations together with the frequency they occurred in our training and testing sets used are given in Table V.

A set of transformations discussed above allows to translate correctly all Polish words apart from a few borrowings, for instance *cyjanek* (cyanide) or foreign names often appear in Polish, for instance a surname *Katz*. However, even translation of such words using only described possibilities does not lead to misunderstanding.

The nature of the text to phoneme transformation issue predestinates it for a neural network approach. To solve it in a classic algorithmic way one would need a set of all relevant transformation rules. Finding a proper set of rules for a given language usually causes a lot of difficulties [12]. For instance, English is particularly difficult to master because of its irregular spelling [6, 13]. Even if a suitable set is given, it is always either huge or incomplete, which leads to impossibility of obtaining 100% accuracy [14] and is therefore no better than a neural net, which never accomplishes 100%. It is always an arduous task to implement it in an algorithmic way. On the

other hand for a given word its phonematic transformation is well known if only a type of pronunciation is given (*eg.* literary or cockney for English). Because training of a neural network requires only a set of examples, it seems to be the convenient method to solve the problem. Finally, the similarities of the learning process of a neural network to that of a human make this approach even more interesting.

The neural network method has also some disadvantages. As it was mentioned, no artificial neural network can give a 100% accuracy. There can be problems with finding a set of training examples that would be statistically suitable, and small enough for the training part to end in a reasonable time. Additionally, finding a proper network architecture is not straightforward and usually a number of trials are needed. The same applies to choosing the right representation of both the input text and the resulting phonemes, one that would at the same time be compact for the net to be of modest size and result in easy input-output mappings. A few neural networks, which we have implemented, are described in the next Section.

3. Implementation of an artificial neural network

We have employed a multilayer feed forward network, which was trained by the backpropagation method (BPM) in a supervised mode [7]. The network consists of layers of neurons — the input layer into which the pattern being recognized is fed, then the hidden layer (or a number of them) which receives a signal from the input layer and propagates it forward, and at the end the output layer which generates the answer. The neurons are fully interconnected on the layer basis, that is every neuron in one layer is connected to each neuron in the next layer, but no neurons in one layer are connected between themselves and there are no connections between neurons in layers that are not neighbouring. During the learning (BPM) phase the generated output is compared to the given correct output (the supervised mode) and an error is computed for each output neuron. Afterwards, the information about this error is “backpropagated” to previous layers in such a way that all neurons receive only a part of the total error, the part they have contributed to. On the basis of this value they adapt themselves in order to lessen the error. The procedure is repeated many times for all examples given in the training set until the total error drops to a value that is small enough. The important part is that the neurons in hidden layers organize themselves in a way as to recognize different features of the input data so that, after training, the network would find these features in inputs that were not presented during training and most probably respond in a correct way.

Such a network would receive a letter at its input and should respond with the correct phoneme. Since the transformation of a single letter into

a phoneme depends on the context (as it was described in Section 2), it is logical to use a sliding window that will analyze simultaneously some number of characters. Then the output should represent the correct phoneme corresponding to the letter in the center of the window.

Now we address the problem of character and phoneme representation. Taking into account the network size and corresponding learning time, it may seem that the representation should be as compact as possible, *e.g.* binary in which each phoneme is represented by a binary number. But then two completely distinct phonemes could become very similar for the network and it would produce erroneous results. We have tried such a representation with only six output neurons, and even less in the hidden layer, but the network could not learn the given examples, not to say about a satisfactory generalization. It turns out that the proper solution is the natural one, *i.e.* the one used in the phonetics as described in the previous Section. Namely each phoneme is represented as a group of some of 14 articulatory features (in Polish language). The abundance of this description is exploited to decrease the probability of confusion of different phonemes.

There is a similar problem with representation of characters. A too compact one is prone to give wrong results. We have found that the network which used a binary representation could not achieve the errors smaller than 50% ! One of the possible solutions is to represent each character as a vector of n neurons with only one activated. This method, which will be referred to as *unary* was also used by Sejnowski [6]. Unfortunately the number of input neurons becomes high ($34 \times 7 = 238$) which lengthens the learning process. Additionally the input is more sensitive to the noise. A change of a value of a single neuron, which corresponds to a noise on a 3% level, results in an input character being equally similar to two different letters and the networks output may be totally wrong. For comparison we have also employed the *distributed* input technique with rather good results.

Each letter is represented with k neurons, where k is much less than n , the number of different characters. Out of 2^k possibilities we take n ¹ different input vectors in_i in such a way that in each one there is approximately $k/2$ active and inactive neurons, and consequently the Hamming distances between all pairs of vectors are similar. This choice has two advantages. First, less neurons are needed to represent a single character — 21 in our case. What follows, is that the number of synapses in the network is lower, thus the number of examples needed to achieve satisfying generalization. On the other hand, it results in a more complicated error function surface and the learning time could be longer. Second, even with a few erroneous

¹ Where n is the number of letters, punctuation, and interword space.

neuron values the resulting configuration still resembles the original one, thus the noise level can be higher. This representation is called *box* by us.

4. Results

We have tested the networks with both the unary and the distributed inputs, and with the output phoneme represented by 14 articulatory features (*velar, high-pitched, low, centralized, vocal, consonantal, nasal, palatal, labial, lateral, fricative, dental, explosive, voiced*) plus 3 neurons needed to represent punctuation marks. The analyzed character was put in the middle of the seven character window. This provided the required information about the context.

The training set consisted of 13000 examples for which the correct phonematic translation was provided. The different cases appeared in the learning set with the frequency corresponding to the commonly used Polish language.

All tested networks had one hidden layer which proved to be sufficient. A network with no hidden layers could not translate text satisfactorily. The network with two hidden layers learned much slower and it had tendencies to "learn by heart", *i.e.* to remember the learning set rather than to generalize.

We have defined two kinds of recognition: one, called *perfect*, when the output contained activations that corresponded to actual phoneme representations; the other, called *correct*, where, if the actual activations did not correspond to any phoneme, than one that was closest, in the sense of the Hamming distance, was chosen, otherwise output consisted of the corresponding phoneme. During the training phase all networks acquired efficiency on the learning set close to 99.0% – 99.9%. The translation efficiency is defined in terms of the ratio of good translations of a character to the total number of transformations performed in the sample. This has changed however when the network was presented with the unknown testing set. The situation is summarized in Table IV. We were studying the performance of a network changing the coding scheme of the input characters (second and third column) and size of the hidden layer. The number of the output neurons was kept fixed at $N_{\text{out}} = 17$ as required by our coding of phonemes by their articulatory features. The fourth column gives the length of the (BPM) training period where one epoch (iteration) is defined as the single pass through the whole learning set. In all cases the training phase was stopped when the preselected accuracy was achieved. The last two columns summarize the performance of the network on the unknown testing set which contained almost 5000 characters. The qualities and speed of learning of different networks are given in Table IV. and in figure 1.

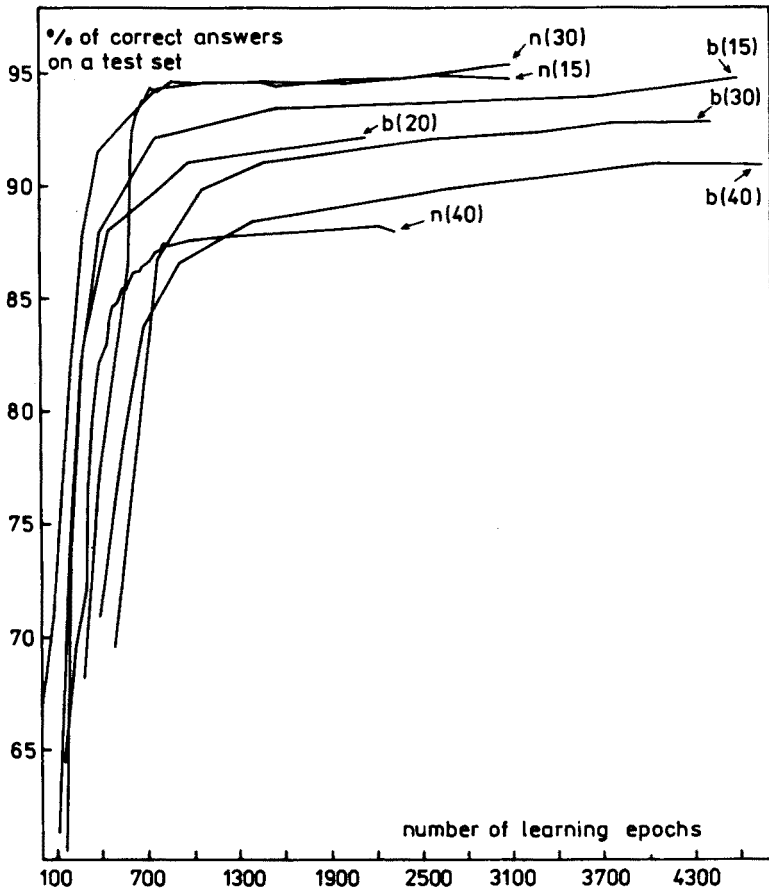


Fig. 1. Learning curves: percent of correct answers versus number of training epochs; $u(n)$ is a unary coded network with n hidden neurons, $b(n)$ is a box coded one.

TABLE IV

The results obtained with different network architectures (different number of hidden neurons)

	type	#input	#hidden	#epochs	time in s	%perfect	%correct
1	unary	238	40	2304	596e3	81.3	88.0
2	unary	238	30	3103	515e3	90.3	95.4
3	unary	238	15	2963	497e3	92.2	94.8
4	box	147	40	4741	898e3	81.6	91.0
5	box	147	35	3989	628e3	83.6	92.2
6	box	147	30	4421	692e3	83.6	92.8
7	box	147	25	2247	395e3	79.5	91.9
8	box	147	20	2108	213e3	86.1	93.6
9	box	147	15	4589	241e3	90.9	94.8
10	box	147	10	3131	258e3	83.0	92.0

TABLE V

Transformations of letters into phonemes in the training set. The number of occurrences is given in brackets.

-	~ (1023)								
a	a (851)								
ą	o (4)	om (3)	on (50)	oń (8)	o ɥ (80)				
b	b (164)	p (31)							
c	ç (17)	č (38)	null (141)	c (109)	ć (55)	č (77)			
ć	č (23)	ć (28)							
d	ç (29)	č (51)	č (17)	c (50)	ć (37)	č (18)	d (203)	t (19)	
e	e (176)								
ę	e (84)	em (11)	en (27)	eń (7)	e ɥ (16)				
f	f (52)	v (25)							
g	g (150)	g (15)	k (61)						
h	x (149)								
i	null (309)	i (432)	i (139)						
j	i (218)								
k	g (34)	k (227)	k (63)						
l	l (262)								
ł	ɥ (226)								
m	m (304)								
n	n (415)	ń (177)	ź (1)						
ń	ń (39)								
o	o (742)								
ó	u (88)								
p	b (40)	p (302)							
r	r (344)	ś (74)	ż (49)						
s	s (268)	ś (68)	ś (167)	z (76)	ż (83)				
ś	ś (113)	ż (82)							
t	d (13)	t (408)							
u	u (192)	ɥ (14)							
w	f (145)	v (355)							
y	i (13)	y (473)							
z	null (322)	s (29)	z (140)	ź (23)					
ź	null (48)	ś (14)	ź (26)						
ż	null (35)	ś (10)	ż (67)						

It can be seen that some networks had 95% of correct translations which is considered as the reasonable threshold for comprehensible spelling. Interestingly, the quality of translation was increasing with the reduction of the size of the hidden layer. This can be also seen in the figure 1. with the learning curves. However below some critical size the performance of the network was decreasing dramatically so that no clear maximum was observed. The maximum efficiency was observed with the unary method of character coding. The distributed input leads to similar results.

The experiments were done at the Technical University of Athens on a Silicon Graphics Power Challenge computer with MIPS R4000 processors. The learning times varied for different architectures (but using the same learning set) from around 59 hours to around 249 hours. However, the satisfactory performances were achieved usually much earlier. The average time needed for the networks to reach, both learning and generalization, errors near their optimal, was around 70 to 80 hours CPU time.

5. Summary

We have studied the neural network implementation of the phonematic translation of Polish text. Although the character — phoneme correspondence is simpler in Polish than in English for example, there exist a number of ambiguities which make application of the neural network suitable and interesting.

Several architectures and coding techniques were tested. Networks with more than one hidden layer were not able to generalize. This property was also observed in other applications [15]. Depending on the number of hidden neurons and on the method of coding of input characters the efficiency of reading the unknown text varied in the 88% to 95% range. In comparison, Sejnowski achieved from 77% to 91% for English language, depending on the network architecture and the training set size [5, 6]. On the other hand, the algorithm described in [14], basing on 329 rules without using any dictionary (for problems occurring with unorthodox events) gained about 96% efficiency, which is comparable to ours. The best choice results with the unary input of 238 neurons and 30 hidden neurons though the network of 147 input neurons and distributed coding performed similarly.

All training curves, as it can be seen in figure 1, are steeply rising within first seven hundred iterations and then slowly saturate during the next few thousand epochs; this feature was also reported in [5]. Therefore it may often suffice in practice to limit the training process to this period of rapid learning.

Many questions remain open. In particular, we would like to study the performance of the distributed networks in the present context. Such

networks were found to give satisfactory results when applied to other problems [16]. The choice of the training set and its effect for the generalization ability of a network requires also further investigations. Finally, it would be very interesting to study in detail the relation between the learning process of the artificial and the natural neural network, *e.g.* by exploiting the existing pedagogical expertise.

Two of us, I. P. and A. B., thank A. Stafylopatis for valuable advises and the hospitality and warm atmosphere that we have enjoyed during our stay at the National Technical University of Athens.

We would like to thank A. Kotański for discussions and for the critical reading of the manuscript.

REFERENCES

- [1] Digital Equipment Corporation, *DECtalk DTC01 Owner's Manual*, Digital Equipment Corporation, Maynard, Mass., document number EK-DTC01-OM-002.
- [2] J. Allen, *From Text to Speech: The MITalk System*, Cambridge University Press, 1985.
- [3] D. Klatt, *J. Acoust. Soc. Am.* **67**, 971 (1980).
- [4] L. Bolc, M. Maksymienko, *Komputerowy system przetwarzania tekstów fonematycznych*, Wydawnictwa Uniwersytetu Warszawskiego, Warszawa 1981.
- [5] T.J. Sejnowski, C.R. Rosenberg, NETtalk: A Parallel Network that Learns to Read Aloud, Johns Hopkins University Department of Electrical Engineering and Computer Science Technical Report 86/01, 1986.
- [6] T.J. Sejnowski, C.R. Rosenberg, *Complex Systems* **1**, 145 1987.
- [7] J. Hertz, A. Krogh, R. Palmer, *Introduction to the Theory of Neural Computation*, Addison-Wesley Publishing Company, 1991.
- [8] R. Jakobson, M. Halle, *Podstawy języka*, PWN, Wrocław 1969, in Polish.
- [9] M. Steffen-Batogowa, *Automatyzacja transkrypcji fonematycznej tekstów polskich*, Warszawa, 1975, in Polish.
- [10] J. Stratyński, *Zarys fonetyki i fonologii*, in Polish.
- [11] M. Karaś, M. Madejowa, *Słownik wymowy polskiej*, Warszawa, Kraków, 1977, in Polish.
- [12] W. Haas, *Phonographic Translation*, Manchester Univ. Press, 1970.
- [13] R.L. Venetky, *The Structure of English Orthography*, The Hague, Mouton, 1970.
- [14] H. Elovitz, R. Johnson, A. McHugh, J. Shore, *IEEE Trans. on Acoustics, Speech, and Signal Processing* **24**, No.6. (1976).

- [15] B. Müller, J. Reinhardt, *Neural Networks An Introduction*, Springer-Verlag, Berlin Heidelberg 1990.
- [16] B. Denby *et al.*, Neural Networks for Triggering, presented at the 1989 IEEE Science Symposium, San Francisco, January 1990; Fermilab report FERMILAB-Conf-90/20.