
A Better Phone Set for the TIMIT Dataset Discovered in Clustering of Listen, Attend and Spell

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Listen, Attend and Spell(LAS)[4] maps a sequence of acoustic spectra directly to a
2 sequence of graphemes, with no explicit internal representation of phones. This
3 paper asks whether LAS can be used as a scientific tool, to discover the phone
4 set of a language whose phone set may be controversial or unknown. Phonemes
5 have a precise linguistic definition, but phones may be defined in any manner
6 that is convenient for speech technology: we propose that a practical phone set
7 is one that can be inferred from speech following certain procedures, but that is
8 also highly predictive of the word sequence. We demonstrate that such a phone
9 set can be inferred by clustering the hidden node activation vectors of an LAS
10 model during training, thus encouraging the model to learn a hidden representation
11 characterized by acoustically compact clusters that are nevertheless predictive of
12 the word sequence. We further define a metric for the quality of a phone set (sum
13 of conditional entropy of the phone set given graphemes, and given acoustics), and
14 demonstrate that according to this metric, the clustered-LAS phone set is better
15 than the original TIMIT[5] phone set.

16 1 Introduction

17 Traditional automatic speech recognition(ASR) usually are composed of multiple components includ-
18 ing an acoustic model, a language model, a pronunciation dictionary, and other possible elements.
19 Recently, modern ASR models implemented based on neural network, such as connectionist temporal
20 classification (CTC)[6] and LAS, accomplished directly speech-to-text with large successes. Since
21 such models generally are not dependent on utilizing specific language models or pronunciation
22 dictionaries, their simple architectures are popular with new researchers trying to enter the speech
23 recognition community. The typical neural-network based models rely on the Recurrent Neural
24 Networks(RNNs), and the key to success of utilizing such deep learning mechanism is to discover
25 hidden representation of the training data.

26 In this work, we take a step further to explore the possibility of defining a new phone set for the
27 TIMIT dataset using the LAS model by incorporating a clustering method to soft align acoustics and
28 graphemes. In LAS model, the Listener takes the input acoustic signals and encodes the signals to a
29 hidden nodes vector, and then the hidden nodes vector feeds into the Speller to generate transcripts.
30 Since the hidden node vector represents the relationship between the words in transcripts and acoustic
31 signals, we cluster these hidden nodes with corresponding trigraphs in the transcripts. In this way, we
32 train the model to learn the underlying relationship between the trigraphs and acoustics.

33 The clustering pairs of hidden nodes and corresponding graphemes are the new defined phone set. We
34 evaluate the new phone set by using an entropy utility function, the sum of the conditional entropy of
35 different contexts given the phone set, and the contexts here refer to both graphemes and acoustics.
36 The experiment reveals that the new phone set discovered by the experiment model better represents
37 the TIMIT dataset under this metric.

38 **2 Related Work**

39 **2.1 Machine learning and attention in ASR**

40 In recent years, machine learning has been extensively researched and applied to many aspects
41 of studies. In the field of ASR, successful models such as CTC, LAS, and [8, 10], and other
42 architectures take advantages of incorporating RNNs. Among these models, CTC and LAS don't
43 require input segmentation and post-processed outputs. CTC generates the labels of sequence of data
44 with RNNs based on the probability distribution given the input sequence during input time steps,
45 whereas LAS neglects input time steps and generates output characters at each output time step using
46 sequence-to-sequence attention mechanism given the transformed hidden nodes vector from input
47 acoustics. Sequence-to-sequence attention mechanism are widely used, and studies such as [9, 2]
48 demonstrate the successes of attention mechanism. These two essential mechanisms greatly benefit
49 current ASR models.

50

51 **2.2 Representations between acoustics and text in speech recognition models**

52 Deep learning works if and only if it's able to find an accurate hidden representation of training data,
53 thereby enabling the system to learn the relationship between the input signal and output words. For
54 conventional ASR, phone recognition and phone segmentation are two important tasks since phone is
55 the smallest temporal unit in speech and serves as an intermediate representation connecting speech
56 and text. Hidden Markov Model(HMM) capture acoustic signal features and decompose vocabulary
57 to context-independent phones [16]. Hybrid HMM-DNN systems use the DNN to compute phone
58 likelihoods, and the HMM to compute phone alignment[11]. Belinkov and Glass [3] investigates the
59 hidden representations of Deep Speech 2[1], and the study shows that the phonetic information loss
60 gradually increases from the bottom layer to the top layer.

61 **3 Model**

62 **3.1 Brief descriptions of LAS model**

63 LAS is an end-to-end speech recognition model that generates the transcripts directly from input
64 acoustic signals without the implementations of multiple submodules of traditional ASR. The basic
65 LAS model includes two modules: a Listener and a Speller. The Listener composes a three layer of
66 pyramidal Bidirectional Long Short Term Memory(pBLSTM)[7], which encodes the input acoustics
67 signals and reduces the input time length to one-eighth of the original. The output of the Listener
68 is represented by hidden nodes vectors \mathbf{h} , then the vectors are fed into the Speller as the input. The
69 Speller is a sequence-to-sequence attention-based LSTM transducer. The attention mechanism of the
70 transducer takes the hidden nodes vectors from the Listener and character distribution from previous
71 step as input and generates a context vector for all attention probabilities of hidden nodes vectors for
72 the current step. The context vector is then used to generate the output character at the current step.

73 **3.2 Experiment LAS model**

74 The Figure 1 shows the overall modified LAS model of the experiment. We are aiming at
75 finding the hidden relationship between the input acoustic signals and output transcripts, so we
76 introduce a clustering component in the original LAS model to encourage the Listener to learn
77 a hidden representation in which frames are grouped into compact clusters. Specifically, for
78 each character correctly inferred by the LSTM transducer from the Speller, we cluster the
79 corresponding maximally attended hidden nodes. The hidden nodes vectors are a cumulative
80 nonlinear transformation of the Mel-Frequency Cepstral Coefficients(MFCCs), and are trained to
81 optimally summarize whatever information about the MFCC is necessary for the Speller to correctly
82 generate output characters. By clustering maximum attended hidden nodes, we force the system to
83 learn groupings of speech frames that have similar hidden node vectors and are also connected to
84 similar output character sequences.

85

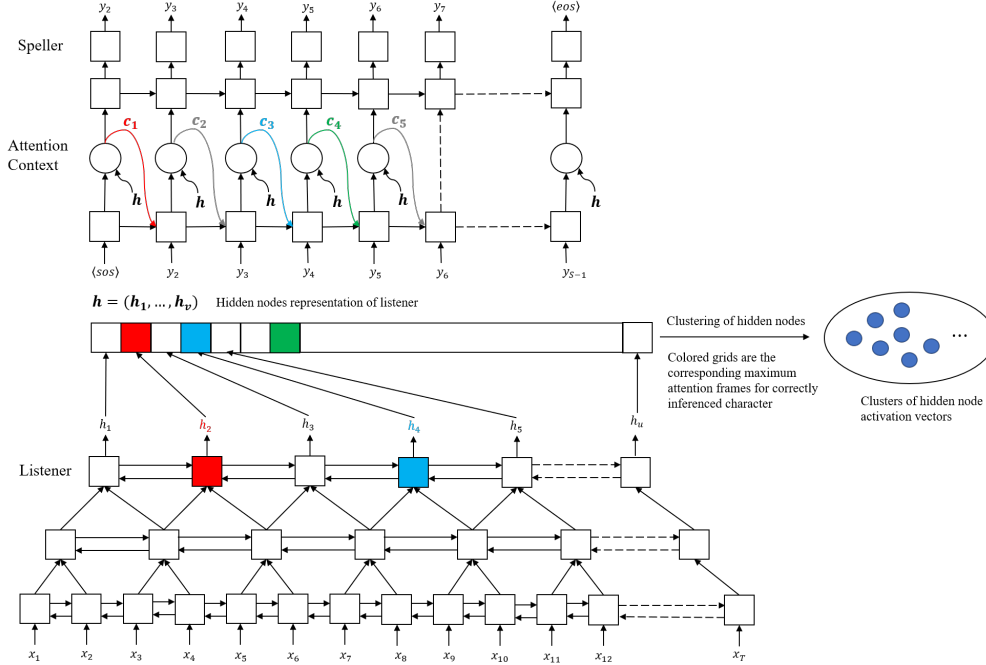


Figure 1: Modified architecture of Listen, attend and spell

The input of the clustering algorithm are the corresponding hidden nodes of the maximum attention frames generated by the AttentionContext vector for correctly inferred character from the Speller. In this figure, $y_2, y_4,$ and y_5 are correctly inferred characters, and their corresponding context vectors of $c_1, c_3,$ and c_4 generate the attention vectors whose maximally attended input frames are $h_2, h_4,$ and h_7 respectively. Thus, $h_2, h_4,$ and h_7 are the input of the cluster.

86 3.3 Learning

87 The modified LAS system can be trained jointly for accurate character output, but also for optimally
 88 clustered internal hidden node vectors. The training criterion of the modified LAS model contains
 89 two parts: word loss and clustering loss. The loss function can be described as the following,

$$\varepsilon = \text{Edit Distance}(y, \hat{y}) + \sum_t \|\mathbf{h}_t - \mu_{k(t)}\|^2$$

90 The first part of the training criterion of the system is the edit distance between the reference transcripts
 91 generated transcripts, and is the error measure used in the standard LAS algorithm; y and \hat{y} refer to
 92 the reference character and generated character respectively. The second part is the squared distance
 93 between each hidden nodes with the centroids of their clusters for the duration of the input sentence.
 94 \mathbf{h}_t is the hidden nodes vector of the input, $\mu_{k(t)}$ is the corresponding cluster of the hidden nodes
 95 vector. By minimizing this error function, we encourage the Listener to learn a hidden embedding,
 96 \mathbf{h}_t , that is useful in predicting the output character y_t , but that can also be clustered into compact
 97 phone-like clusters with centroids μ_k .

98 3.4 Clustering method

99 The clustering method in the modified LAS model has almost the same mechanism of k-means
 100 clustering except that the input varies after every step since the hidden nodes change for every batch
 101 during training stage. The objective of the clustering method is to minimize the clustering loss from
 102 the loss function.

103 The centroids of the clusters are randomly initialized with normal distribution. The hidden nodes
 104 are clustered and labeled for a certain number of iterations for every epoch. Then the centroids are
 105 updated and kept for the next epoch. After every epoch, the empty clusters that are never labeled
 106 for the past epoch will be deleted and replaced by splitting the largest labeled clusters by scaling the
 107 original centroids 0.01 and 0.99 of the original clustering centroids.

108 4 Experiment

109 4.1 Dataset descriptions

110 Two datasets are used to perform the experiment. English speech recognition training corpus of
111 TED-LIUMv2(TEDLIUM)[15] is used to pre-train the LAS model. The TEDLIUM dataset was
112 made from audio talks and transcripts from TED website. There are 1495 audio talks with aligned
113 transcripts in the dataset.

114 The TIMIT dataset is used to train the experiment LAS model. The TIMIT dataset comes with its
115 self-defined dictionary and phoneme alignment transcripts for the audio talks. During training, audio
116 and transcripts of one female and one male are selected from each dialect region of the test dataset as
117 development set, and the rest of the test dataset remains as test set. A new phone set is discovered for
118 the TIMIT dataset and compared with the reference phone set.

119 4.2 Preprocessing of dataset

120 Both dataset are preprocessed using MFCCs algorithm. The raw acoustic signals in the dataset are
121 framed by 10ms each, and the sampling rate of the input signals is 16000 Hz. The power spectrum is
122 calculated for each frame by using periodogram estimate, the squared magnitude of Discrete Fourier
123 Transform(DFT) of original acoustic signals. Forty filters in Mel-spaced filterbank are applied to the
124 power spectrum, and log filterbank energies are computed by taking the log of the power spectrum.
125 Then Discrete Cosine Transform(DCT) of these forty log filterbank energies give output of the
126 cepstrum coefficients.

127 4.3 Experimental settings

128 The implementation of the basic LAS model is based on the toolkit eXtensible Neural Machine
129 Translation(XNMT)[13] using Dynet framework[12]. The learning rate of the Adam optimizer is
130 initialized to 0.01 and reduced to half of the original learning rate if WER of development set isn't
131 improved after 3 epochs. Other parameters are consistent as indicated by the original paper. The
132 hidden dimension of pLSTM is 512, which is the dimension of the hidden nodes vector. The Attender
133 has hidden dimension 128. The dropout rate of the entire neural network is 0.3. The Speller uses a
134 beam search with size 20 is used to infer test transcripts.

135 The experiment model is modified by introducing a new clustering module. The LAS model has
136 been pre-trained for about 300 epochs. Starting with the pre-trained model, the experiment model is
137 then trained for 100 epochs; the learning rate of the Adam optimizer remains at 0.03. The number of
138 iterations for each clustering step is 20, and the dimension of each cluster centroid is the same as
139 the dimension of the hidden nodes vector, which is 512 in this case. The number of clusters is 100,
140 which is roughly twice as the number of English phonemes.

141 5 Results and discussions

142 5.1 Error measurements of experiment LAS model

143 Upon convergence, the pre-trained model of LAS has word error rate(WER) 16.72% and character
144 error rate(CER) 8.46% on the test dataset. With the pre-trained LAS model, the experiment model
145 has the final WER 26.99% and CER 10.67%. By the training criterion, clustering loss is 0.359 and
146 maximum likelihood estimation loss is 0.948.

147

148 5.2 Comparisons of new discovered phone set and reference phone set

149 For the experimental LAS model, 100 clusters are used to discover a new phone set for TIMIT. For
150 all generated transcripts of test set, every character in the transcripts is assigned to a cluster using
151 the k-means cluster criterion, and the closest μ_k of the corresponding hidden node \mathbf{h}_t is the one that
152 each character assigned to. The top five most frequently assigned trigraphs for each cluster vote to
153 determine the phone label of each cluster from phones in the TIMIT dataset. There are certainly some
154 ambiguous cases when we try to identify the phones for one cluster. For example, the top five most
155 frequently labeled trigraphs for one cluster are " wh", "ere", "whe", "wer", and " we". The cluster

Table 1: List of both phone sets discovered by experiment model and reference dictionary

phone categories	phone set of experiment model	phone set of reference dictionary
stops	b, d, g, k, p, t	b, d, g, k, p, t
affricates	ch, jh	ch, jh
fricatives	f, s, sh, th, v, z	dh , f, s, sh, th, v, z, zh
nasals	en, m, n, ng	em , en, eng , m, n, ng
semivowels and glides	el, hh, hv , r, w, l	el, hh, l, r, w, y
vowels	ae, ao, ax, axr, aw, ay, eh, er, ey, ih, ix, iy, ow, oy, uw	aa , ae, ah , ao, aw, ax, axr, ay, eh er, ey, ih, ix, iy, ow, oy, uh , uw
non-speech event	h#	h#

Table 2: Entropy of the distribution $P(\text{phones}|\text{graphemes})$ and $P(\text{phones}|\text{acoustics})$ for both experiment and reference phone sets

System	$H(\text{phones} \text{graphemes})$	$H(\text{phones} \text{acoustics})$
Experiment	0.0212	0.00117
Reference	0.0242	0.00136

156 certainly captures the similar pronunciations of the word "where", but the word may not be described
 157 using one single phoneme. In such cases, phone labels of clusters are edited by hand. For this cluster,
 158 we assign the phone label as "w". The final phone set discovered by the clusters of the experimental
 159 system contain 40 unique phones.

160 Since the reference transcripts of the TIMIT dataset contain the actual pronunciations of the phones,
 161 the phones of the transcripts are very different from the ones used in the TIMIT dictionary. In order
 162 to find a phone set that represents the reference transcripts, we utilize the function phonetisaurus-
 163 align in toolkit Phonetisaurus G2P[14] to generate the alignment between each character and the
 164 corresponding phone of the reference transcript. The stress markers of the TIMIT dictionary are
 165 eliminated. The reference phone set contain 46 unique phones.

166 The experiment and reference phone sets are as shown in the table1.

167

168 5.3 Entropy measurement

169

170 Entropy is commonly used to measure the randomness or disorder of a system. The output of the
 171 experiment model is evaluated by calculating the conditional entropy given different contexts for
 172 experiment and reference set of phones. The contexts include both graphemes and acoustics. The
 173 graphemes constitute of all possibilities of trigraphs in English, and the acoustics include all unique
 174 frames in the test dataset from TIMIT corpus. The conditional entropy is calculated as the following,

$$H(\text{phones}|\text{contexts}) = - \sum_{x \in \text{contexts}} p(x) \sum_{y \in \text{phones}} p(y|x) \log p(y|x)$$

175 Specifically, for calculating $H(\text{phones}|\text{graphemes})$, all possibilities of trigraphs are considered as the
 176 contexts of graphemes. The possibilities of trigraphs are calculated as all of length three permutations
 177 of 26 English letters and special tokens "" and "-" appeared in both generated and reference
 178 transcripts.

$$p(x) = \frac{\text{number of occurrences for trigraph } x \text{ in transcripts}}{\text{total number of trigraphs in transcripts}}$$

179

$$p(y|x) = \frac{\text{number of occurrences for phone } y \text{ given trigraph } x \text{ in transcripts} + k}{\text{total number of trigraphs in transcripts} + k \times \text{number of phones in defined phone set}}$$

180 Similarly, for calculating $H(\text{phones}|\text{acoustics})$, unique acoustic frame is sorted out by measuring
 181 squared differences of all acoustic frames among each other in test dataset. Within tolerance of
 182 1, similar acoustic frames are treated as the same frame in calculation. Thus, the prior acoustic
 183 distribution can be approximated by

$$p(x) \approx \frac{\text{number of occurrences of acoustic frame } x}{\text{total number of acoustic frames in test set}}$$

$$p(y|x) = \frac{\text{number of occurrences for phone } y \text{ given acoustic frame } x + k}{\text{total number of acoustic frames in test data} + k \times \text{number of phones in defined phone set}}$$

185 Laplace smoothing is applied for both conditional entropy calculations with the smoothing factor
 186 $k = 1$.

187 From Table2, the sum of the conditional entropy given both contexts of experiment phone set is less
 188 than that of the reference phone set, since $0.0212 + 0.00117 = 0.02237 < 0.0242 + 0.00136 =$
 189 0.02556 .

190 6 Conclusions

191 We defined a new phone set for the TIMIT dataset based on incorporating the clustering mechanism
 192 into the original LAS model. The learning criterion for the experiment model is composed of two parts:
 193 edit distance between generated transcripts and reference transcripts and squared distance between
 194 clustered hidden nodes and corresponding centroids of clusters. The learning criterion balances the
 195 learning objectives of the system – reducing the WER of generated transcripts meanwhile grouping
 196 the hidden vectors into compact clusters. The model is pre-trained with a larger dataset, TEDLIUM,
 197 and then trained on the TIMIT dataset for the experiment. The experiment result is evaluated by
 198 defining a utility function, the sum of the conditional entropy of graphemes given phones and the
 199 conditional entropy of acoustics given phones. We showed that the experiment phone set has both
 200 lower "grapheme entropy" and "acoustic entropy". Thus, we can claim that the phone set discovered
 201 by the experiment is better than the reference phone set in the TIMIT dataset based on this criterion.

202 Acknowledgments

203 We thank for everyone who provides insightful discussions with us for conducting research.

204 References

- 205 [1] D. Amodei, S. Ananthanarayanan, R. Anubhai, J. Bai, E. Battenberg, C. Case, J. Casper,
 206 B. Catanzaro, Q. Cheng, G. Chen, J. Chen, J. Chen, M. Chen, Z. and Chrzanowski, A. Coates,
 207 G. Diamos, K. Ding, N. Du, E. Elsen, J. Engel, W. Fang, L. Fan, C. Fougner, L. Gao, C. Gong,
 208 A. Hannun, T. Han, L. Johannes, B. Jiang, C. Ju, B. Jun, P. Legresley, L. Lin, J. Liu, Y. Liu,
 209 W. Li, X. Li, D. Ma, S. Narang, A. Ng, S. Ozair, Y. Peng, R. Prenger, S. Qian, Z. Quan,
 210 and Raima. Deep speech 2: End-to-end speech recognition in English and Mandarin. In
 211 *33rd International Conference on Machine Learning, ICML 2016*, number 33rd International
 212 Conference on Machine Learning, ICML 2016, pages 312–321, (1)Baidu Silicon Valley AI Lab,
 213 2016.
- 214 [2] D. Bahdanau, J. Chorowski, D. Serdyuk, P. Brakel, and Y. Bengio. End-to-end attention-based
 215 large vocabulary speech recognition. In *2016 IEEE International Conference on Acoustics,
 216 Speech and Signal Processing (ICASSP)*, pages 4945–4949, Mar 2016.
- 217 [3] Y. Belinkov and J. Glass. Analyzing hidden representations in end-to-end automatic speech
 218 recognition systems. In *Proceedings of Advances in Neural Information Processing Systems 30
 219 (NIPS 2017)*, 2017.
- 220 [4] W. Chan, N. Jaitly, Q. V. Le, and O. Vinyals. Listen, attend and spell: A neural network for
 221 large vocabulary conversational speech recognition. In *2016 IEEE International Conference On
 222 Acoustics, Speech And Signal Processing (ICASSP)*, pages 4960–4964, May 2016.
- 223 [5] J. Garofolo, L. Lamel, W. Fisher, D. P. Jonathan Fiscus, N. Dahlgren, and V. Zue. TIMIT
 224 acoustic-phonetic continuous speech corpus LDC93S1, 1993. Web Download.

- 225 [6] A. Graves, S. Fernández, and F. Gomez. Connectionist temporal classification: Labelling
226 unsegmented sequence data with recurrent neural networks. In *Proceedings of the International*
227 *Conference on Machine Learning, ICML 2006*, pages 369–376, 2006.
- 228 [7] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 9(8):
229 1735–1780, 1997.
- 230 [8] T. Hori, S. Watanabe, Y. Zhang, and W. Chan. Analyzing hidden representations in end-to-
231 end automatic speech recognition systems. In *Proceedings of the Annual Conference of the*
232 *International Speech Communication Association, INTERSPEECH, 2017*, pages 949–953,
233 2017.
- 234 [9] S. Kim, T. Hori, and S. Watanabe. Joint ctc-attention based end-to-end speech recognition
235 using multi-task learning. In *ICASSP, IEEE International Conference on Acoustics, Speech and*
236 *Signal Processing - Proceedings*, pages 4835–4839, 2017.
- 237 [10] Y. Miao, M. Gowayyed, and F. Metze. Eesen: End-to-end speech recognition using deep RNN
238 models and wfst-based decoding. In *2015 IEEE Workshop on Automatic Speech Recognition*
239 *and Understanding (ASRU)*, pages 167–174, Dec 2015.
- 240 [11] N. Morgan and H. Bourlard. Continuous speech recognition. In *IEEE Signal Processing*
241 *Magazine IEEE Signal Process. Mag. Signal Processing Magazine, IEEE.*, volume 12, May
242 1995.
- 243 [12] G. Neubig, C. Dyer, Y. Goldberg, A. Matthews, W. Ammar, A. Anastasopoulos, M. Balles-
244 teros, D. Chiang, D. Clothiaux, T. Cohn, K. Duh, M. Faruqui, C. Gan, D. Garrette, Y. Ji,
245 L. Kong, A. Kuncoro, G. Kumar, C. Malaviya, P. Michel, Y. Oda, M. Richardson, N. Saphra,
246 S. Swayamdipta, and P. Yin. Dynet: The dynamic neural network toolkit. *arXiv preprint*
247 *arXiv:1701.03980*, 2017.
- 248 [13] G. Neubig, M. Sperber, X. Wang, M. Felix, A. Matthews, S. Padmanabhan, Y. Qi, D. S. Sachan,
249 P. Arthur, P. Godard, J. Hewitt, R. Riad, and L. Wang. XNMT: The extensible neural machine
250 translation toolkit. In *Conference of the Association for Machine Translation in the Americas*
251 *(AMTA) Open Source Software Showcase*, Boston, March 2018.
- 252 [14] J. Novak, P. Dixon, N. Minematsu, K. Hirose, C. Hori, and H. Kashioka. Improving WFST-
253 based G2P conversion with alignment constraints and RNNLM n-best rescoring. In *Interspeech*,
254 2012.
- 255 [15] A. Rousseau, P. Deléglise, and Y. Estève. Enhancing the TED-LIUM corpus with selected
256 data for language modeling and more TED talks. In *Proceedings of the Ninth International*
257 *Conference on Language Resources and Evaluation (LREC'14)*, May 2014.
- 258 [16] T. Vaich and A. Cohen. HMM phoneme recognition with supervised training and viterbi
259 algorithm. In *Eighteenth Convention of Electrical and Electronics Engineers in Israel Electrical*
260 *and Electronics Engineers in Israel, 1995.*, Mar 1995.