# A DENSITY RATIO APPROACH TO LANGUAGE MODEL FUSION IN END-TO-END AUTOMATIC SPEECH RECOGNITION

*Erik McDermott, Hasim Sak, Ehsan Variani*

Google Inc., USA

{erikmcd,hasim,variani}@google.com

## ABSTRACT

This article describes a density ratio approach to integrating external Language Models (LMs) into end-to-end models for Automatic Speech Recognition (ASR). Applied to a Recurrent Neural Network Transducer (RNN-T) ASR model trained on a given domain, a matched in-domain RNN-LM, and a target domain RNN-LM, the proposed method uses Bayes' Rule to define RNN-T posteriors for the target domain, in a manner directly analogous to the classic hybrid model for ASR based on Deep Neural Networks (DNNs) or LSTMs in the Hidden Markov Model (HMM) framework (Bourlard & Morgan, 1994). The proposed approach is evaluated in cross-domain and limited-data scenarios, for which a significant amount of target domain text data is used for LM training, but only limited (or no) {audio, transcript} training data pairs are used to train the RNN-T. Specifically, an RNN-T model trained on paired audio & transcript data from YouTube is evaluated for its ability to generalize to Voice Search data. The Density Ratio method was found to consistently outperform the dominant approach to LM and end-to-end ASR integration, Shallow Fusion.

***Index Terms***— End-to-end models, Automatic Speech Recognition

## 1. INTRODUCTION

End-to-end models such as Listen, Attend & Spell (LAS) [1] or the Recurrent Neural Network Transducer (RNN-T) [2] are sequence models that directly define $P(W|X)$, the posterior probability of the word or subword sequence $W$ given an audio frame sequence $X$, with no chaining of sub-module probabilities. State-of-the-art, or near state-of-the-art results have been reported for these models on challenging tasks [3, 4].

End-to-end ASR models in essence do not include independently trained symbols-only or acoustics-only sub-components. As such, they do not provide a clear role for language models $P(W)$ trained only on text/transcript data. There are, however, many situations where we would like to use a separate LM to complement or modify a given ASR system. In particular, no matter how plentiful the paired {audio, transcript} training data, there are typically orders of
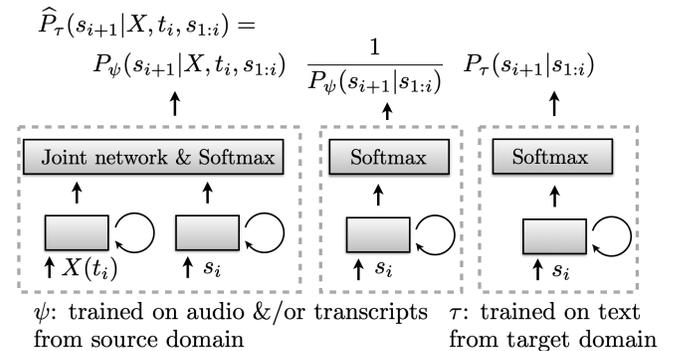
$$\widehat{P}_\tau(s_{i+1}|X, t_i, s_{1:i}) =$$
$$P_\psi(s_{i+1}|X, t_i, s_{1:i}) \quad \frac{1}{P_\psi(s_{i+1}|s_{1:i})} \quad P_\tau(s_{i+1}|s_{1:i})$$



$\psi$: trained on audio &/or transcripts from source domain   $\tau$: trained on text from target domain

**Fig. 1**. Estimating a target domain pseudo-posterior via combination of source domain RNN-T, source domain RNN-LM, and target domain RNN-LM.

magnitude more text-only data available. There are also many practical applications of ASR where we wish to adapt the language model, e.g., biasing the recognition grammar towards a list of specific words or phrases for a specific context.

The research community has been keenly aware of the importance of this issue, and has responded with a number of approaches, under the rubric of "Fusion". The most popular of these is "Shallow Fusion" [5, 6, 7, 8, 9], which is simple log-linear interpolation between the scores from the end-to-end model and the separately-trained LM. More structured approaches, "Deep Fusion" [10], "Cold Fusion" [11] and "Component Fusion" [12] jointly train an end-to-end model with a pre-trained LM, with the goal of learning the optimal combination of the two, aided by gating mechanisms applied to the set of joint scores. These methods have not replaced the simple Shallow Fusion method as the go-to method in most of the ASR community. Part of the appeal of Shallow Fusion is that it does not require model retraining – it can be applied purely at decoding time. The Density Ratio approach proposed here can be seen as an extension of Shallow Fusion, sharing some of its simplicity and practicality, but offering a theoretical grounding in Bayes' rule.

After describing the historical context, theory and practical implementation of the proposed Density Ratio method, this article describes experiments comparing the method to

Shallow Fusion in a cross-domain scenario. An RNN-T model was trained on large-scale speech data with semi-supervised transcripts from YouTube videos, and then evaluated on data from a live Voice Search service, using an RNN-LM trained on Voice Search transcripts to try to boost performance. Then, exploring the transition between cross-domain and in-domain, limited amounts of Voice Search speech data were used to fine-tune the YouTube-trained RNN-T model, followed by LM fusion via both the Density Ratio method and Shallow Fusion. The ratio method was found to produce consistent gains over Shallow Fusion in all scenarios examined.

## 2. A BRIEF HISTORY OF LANGUAGE MODEL INCORPORATION IN ASR

**Generative models and Bayes' rule**. The Noisy Channel Model underlying the origins of statistical ASR [13] used Bayes' rule to combine generative models of both the acoustics $p(X|W)$ and the symbol sequence $P(W)$:

$$p(X|W) = \sum_{\mathbf{s} \in S_W} p(X|\mathbf{s})p(\mathbf{s}|W) = \sum_{\mathbf{s} \in S_W} \prod_t p(\mathbf{x}_t|\mathbf{s}(t))$$
$$P(W|X) = p(X|W)P(W)/p(X) \qquad (1)$$

for an acoustic feature vector sequence $X = \mathbf{x}_1, ..., \mathbf{x}_T$ and a word or sub-word sequence $W = s_1, ..., s_U$ with possible time alignments $S_W = \{..., \mathbf{s}, ...\}$. ASR decoding then uses the posterior probability $P(W|X)$. A prior $p(\mathbf{s}|W)$ on alignments can be implemented e.g. via a simple 1st-order state transition model. Though lacking in discriminative power, the paradigm provides a clear theoretical framework for decoupling the acoustic model (AM) $p(X|W)$ and LM $P(W)$.

**Hybrid model for DNNs/LSTMs within original ASR framework**. The advent of highly discriminative Deep Neural Networks (DNNs) [14, 15, 16, 17, 18] and Long Short Term Memory models (LSTMs) [19, 20] posed a challenge to the original Noisy Channel Model, as they produce phoneme- or state- level posteriors $P(\mathbf{s}(t)|\mathbf{x}_t)$, not acoustic likelihoods $p(\mathbf{x}_t|\mathbf{s}(t))$. The "hybrid" model [21] proposed the use of scaled likelihoods, i.e. posteriors divided by separately estimated state priors $P(w)$. For bidirectional LSTMs, the scaled-likelihood over a particular alignment $\mathbf{s}$ is taken to be

$$P'(X|\mathbf{s}) \equiv k(X) \prod_t P(\mathbf{s}(t)|X)/P(\mathbf{s}(t)), \qquad (2)$$

using $k(X)$ to represent a $p(X)$-dependent term shared by all hypotheses $W$, that does not affect decoding. This "pseudo-generative" score can then be plugged into the original model of Eq. (1) and used for ASR decoding with an arbitrary LM $P(W)$. For much of the ASR community, this approach still constitutes the state-of-the-art [3, 22, 23].

**Shallow Fusion**. The most popular approach to LM incorporation for end-to-end ASR is a linear interpolation,

$$\text{Score}(X, W) = \log P(W|X) + \lambda \log P(W) + \beta|W|, \qquad (3)$$

with no claim to direct interpretability according to probability theory, and often a reward for sequence length $|W|$, scaled by a factor $\beta$ [6, 8, 9, 24].

## 3. LANGUAGE MODEL INCORPORATION INTO END-TO-END ASR, USING BAYES' RULE

### 3.1. A Sequence-level Hybrid Pseudo-Generative Model

The model makes the following assumptions:

1. The source domain $\psi$ has some true joint distribution $P_\psi(W, X)$ over text and audio;

2. The target domain $\tau$ has some other true joint distribution $P_\tau(W, X)$;

3. A source domain end-to-end model (e.g. RNN-T) captures $P_\psi(W|X)$ reasonably well;

4. Separately trained LMs (e.g. RNN-LMs) capture $P_\psi(W)$ and $P_\tau(W)$ reasonably well;

5. $p_\psi(X|W)$ is roughly equal to $p_\tau(X|W)$, i.e. the two domains are acoustically consistent; and

6. The target domain posterior, $P_\tau(W|X)$, is unknown.

The starting point for the proposed **Density Ratio Method** is then to express a "hybrid" scaled acoustic likelihood for the source domain, in a manner paralleling the original hybrid model [21]:

$$p_\psi(X|W) = p_\psi(X)P_\psi(W|X)/P_\psi(W). \qquad (4)$$

Similarly, for the target domain:

$$p_\tau(X|W) = p_\tau(X)P_\tau(W|X)/P_\tau(W). \qquad (5)$$

Given the stated assumptions, one can then estimate the target domain posterior as:

$$\widehat{P}_\tau(W|X) = k(X) \frac{P_\tau(W)}{P_\psi(W)} P_\psi(W|X), \qquad (6)$$

with $k(X) = p_\psi(X)/p_\tau(X)$ shared by all hypotheses $W$, and the ratio $P_\tau(W)/P_\psi(W)$ (really a probablity mass ratio) giving the proposed method its name.

In essence, this model is just an application of Bayes' rule to end-to-end models and separate LMs. The approach can be viewed as the sequence-level version of the classic hybrid model [21]. Similar use of Bayes' rule to combine ASR scores with RNN-LMs has been described elsewhere, e.g. in work connecting grapheme-level outputs with word-level LMs [7, 25, 26]. However, to our knowledge this approach has not been applied to end-to-end models in cross-domain settings, where one wishes to leverage a language model from the target domain. For a perspective on a "pure" (non-hybrid) deep generative approach to ASR, see [27].

## 3.2. Top-down fundamentals of RNN-T

The *RNN Transducer (RNN-T)* [2] defines a sequence-level posterior $P(W|X)$ for a given acoustic feature vector sequence $X = \mathbf{x}_1, ..., \mathbf{x}_T$ and a given word or sub-word sequence $W = s_1, ..., s_U$ in terms of possible alignments $S_W = \{..., (\mathbf{s}, \mathbf{t}), ...\}$ of $W$ to $X$. The tuple $(\mathbf{s}, \mathbf{t})$ denotes a specific alignment sequence, a symbol sequence and corresponding sequence of time indices, consistent with the sequence $W$ and utterance $X$. The symbols in $\mathbf{s}$ are elements of an expanded symbol space that includes optional, repeatable blank symbols used to represent acoustics-only path extensions, where the time index is incremented, but no non-blank symbols are added. Conversely, non-blank symbols are only added to a partial path time-synchronously. (I.e., using $i$ to index elements of $\mathbf{s}$ and $\mathbf{t}$, $t_{i+1} = t_i + 1$ if $s_{i+1}$ is blank, and $t_{i+1} = t_i$ if $s_{i+1}$ is non-blank). $P(W|X)$ is defined by summing over alignment posteriors:

$$P(W|X) = \sum_{(\mathbf{s},\mathbf{t}) \in S_W} P(\mathbf{s}, \mathbf{t}|X) \qquad (7)$$

$$P(\mathbf{s}, \mathbf{t}|X) = \prod_i P(s_{i+1}|X, t_i, s_{1:i}). \qquad (8)$$

Finally, $P(s_{i+1}|X, t_i, s_{1:i})$ is defined using an LSTM-based acoustic encoder with input $X$, an LSTM-based label encoder with non-blank inputs $s$, and a feed-forward joint network combining outputs from the two encoders to produce predictions for all symbols $s$, including the blank symbol.

The Forward-Backward algorithm can be used to calculate Eq. (7) efficiently during training, and Viterbi-based beam search (based on the argmax over possible alignments) can be used for decoding when $W$ is unknown [2, 28].

## 3.3. Application of Shallow Fusion to RNN-T

Shallow Fusion (Eq. (3)) can be implemented in RNN-T for each time-synchronous non-blank symbol path extension. The LM score corresponding to the same symbol extension can be "fused" into the log-domain score used for decoding:

$$\text{Score}(s_{i+1}|X, t_i, s_{1:i}) = \log P(s_{i+1}|X, t_i, s_{1:i})$$
$$+ \lambda \log P(s_{i+1}|s_{1:i}) + \beta. \qquad (9)$$

This is only done when the hypothesized path extension $s_{i+1}$ is a non-blank symbol; the decoding score for blank symbol path extensions is the unmodified $\log P(s_{i+1}|X, t_i, s_{1:i})$.

## 3.4. Application of the Density Ratio Method to RNN-T

Eq. (6) can be implemented via an estimated RNN-T "pseudo-posterior", when $s_{i+1}$ is a non-blank symbol:

$$\widehat{P}_\tau(s_{i+1}|X, t_i, s_{1:i}) = \frac{P_\tau(s_{i+1}|s_{1:i})}{P_\psi(s_{i+1}|s_{1:i})} P_\psi(s_{i+1}|X, t_i, s_{1:i}). \qquad (10)$$

This estimate is not normalized over symbol outputs, but it plugs into Eq. (8) and Eq. (7) to implement the RNN-T version of Eq. (6). In practice, scaling factors $\lambda_\psi$ and $\lambda_\tau$ on the LM scores, and a non-blank reward $\beta$, are used in the final decoding score:

$$\text{Score}(s_{i+1}|X, t_i, s_{1:i}) = \log P_\psi(s_{i+1}|X, t_i, s_{1:i})$$
$$+\lambda_\tau \log P_\tau(s_{i+1}|s_{1:i}) - \lambda_\psi \log P_\psi(s_{i+1}|s_{1:i}) + \beta. \qquad (11)$$

## 3.5. Implementation

The ratio method is very simple to implement. The procedure is essentially to:

1. Train an end-to-end model such as RNN-T on a given source domain training set $\psi$ (paired audio/transcript data);

2. Train a neural LM such as RNN-LM on text transcripts from the same training set $\psi$;

3. Train a second RNN-LM on the target domain $\tau$;

4. When decoding on the target domain, modify the RNN-T output by the ratio of target/training RNN-LMs, as defined in Eq. (11), and illustrated in Fig. 1.

The method is purely a decode-time method; no joint training is involved, but it does require tuning of the LM scaling factor(s) (as does Shallow Fusion). A held-out set can be used for that purpose.

## 4. TRAINING, DEVELOPMENT AND EVALUATION DATA

### 4.1. Training data

The following data sources were used to train the RNN-T and associated RNN-LMs in this study.

**Source-domain baseline RNN-T**: approximately 120M segmented utterances (190,000 hours of audio) from YouTube videos, with associated transcripts obtained from semi-supervised caption filtering [29].

**Source-domain normalizing RNN-LM**: transcripts from the same 120M utterance YouTube training set. This corresponds to about 3B tokens of the sub-word units used (see below, Section 5.1).

**Target-domain RNN-LM**: 21M text-only utterance-level transcripts from anonymized, manually transcribed audio data, representative of data from a Voice Search service. This corresponds to about 275M sub-word tokens.

**Target-domain RNN-T fine-tuning data**: 10K, 100K, 1M and 21M utterance-level {audio, transcript} pairs taken from anonymized, transcribed Voice Search data. These fine-tuning sets roughly correspond to 10 hours, 100 hours, 1000 hours and 21,000 hours of audio, respectively.

**Table 1**. Training set size and test set perplexity for the morph-level RNN-LMs (training domain → testing domain) used in this study.

| Model | # Tr. tokens | Test PPL |
|---|---|---|
| YouTube → YouTube | 2.98B | 8.94 |
| YouTube → Voice Search | 2.98B | 36.5 |
| Voice Search → Voice Search | 275M | 11.1 |

### 4.2. Dev and Eval Sets

The following data sources were used to choose scaling factors and/or evaluate the final model performance.

**Source-domain Eval Set (YouTube)**. The in-domain performance of the YouTube-trained RNN-T baseline was measured on speech data taken from Preferred Channels on YouTube [30]. The test set is taken from 296 videos from 13 categories, with each video averaging 5 minutes in length, corresponding to 25 hours of audio and 250,000 word tokens in total.

**Target-domain Dev & Eval sets (Voice Search)**. The Voice Search dev and eval sets each consist of approximately 7,500 anonymized utterances (about 33,000 words and corresponding to about 8 hours of audio), distinct from the fine-tuning data described earlier, but representative of the same Voice Search service.

### 5. CROSS-DOMAIN EVALUATION: YOUTUBE-TRAINED RNN-T → VOICE SEARCH

The first set of experiments uses an RNN-T model trained on {audio, transcript} pairs taken from segmented YouTube videos, and evaluates the cross-domain generalization of this model to test utterances taken from a Voice Search dataset, with and without fusion to an external LM.

### 5.1. RNN-T and RNN-LM model settings

The overall structure of the models used here is as follows:

**RNN-T**:

- Acoustic features: 768-dimensional feature vectors obtained from 3 stacked 256-dimensional logmel feature vectors, extracted every 20 msec from 16 kHz waveforms, and sub-sampled with a stride of 3, for an effective final feature vector step size of 60 msec.

- Acoustic encoder: 6 LSTM layers x (2048 units with 1024-dimensional projection); bidirectional.

- Label encoder (aka "decoder" in end-to-end ASR jargon): 1 LSTM layer x (2048 units with 1024-dimensional projection).

| WER | Sequence length scaling factor ($\beta$) | | | |
|---|---|---|---|---|
| $\lambda$ | 0.40 | 0.50 | 0.60 | 0.70 |
| 0.15 | 15.84 | 16.53 | 17.21 | 17.84 |
| 0.20 | 15.05 | 15.33 | 15.80 | 16.40 |
| 0.25 | 14.83 | 14.88 | 15.14 | 15.38 |
| 0.30 | 14.97 | 14.84 | 14.83 | 14.96 |
| 0.35 | 15.07 | 15.05 | 14.92 | 14.87 |
| 0.40 | 15.84 | 15.42 | 15.34 | 15.27 |
| 0.45 | 16.81 | 16.56 | 16.25 | 15.96 |

**Fig. 2**. Dev set WERs for Shallow Fusion LM scaling factor $\lambda$ vs. sequence length scaling factor $\beta$.

- RNN-T joint network hidden dimension size: 1024.

- Output classes: 10,000 sub-word "morph" units [31] , input via a 512-dimensional embedding.

- Total number of parameters: approximately 340M

**RNN-LMs** for both source and target domains were set to match the RNN-T decoder structure and size:

- 1 layer x (2048 units with 1024-dimensional projection).

- Output classes: 10,000 morphs (same as the RNN-T).

- Total number of parameters: approximately 30M.

The RNN-T and the RNN-LMs were independently trained on 128-core tensor processing units (TPUs) using full unrolling and an effective batch size of 4096. All models were trained using the Adam optimization method [32] for 100K-125K steps, corresponding to about 4 passes over the 120M utterance YouTube training set, and 20 passes over the 21M utterance Voice Search training set. The trained RNN-LM perplexities (shown in Table 1) show the benefit to Voice Search test perplexity of training on Voice Search transcripts.

### 5.2. Experiments and results

In the first set of experiments, the constraint $\lambda_\psi = \lambda_\tau$ was used to simplify the search for the LM scaling factor in Eq. 11. Fig. 2 and Fig. 3 illustrate the different relative sensitivities of WER to the LM scaling factor(s) for Shallow Fusion and the Density Ratio method, as well as the effect of the RNN-T sequence length scaling factor, measured on the dev set.

The LM scaling factor affects the relative value of the symbols-only LM score vs. that of the acoustics-aware RNN-T score. This typically alters the balance of insertion vs. deletion errors. In turn, this effect can be offset (or amplified) by the sequence length scaling factor $\beta$ in Eq. (3), in the case

| WER | Sequence length scaling factor (β) | | | |
|---|---|---|---|---|
| λ | -0.40 | -0.30 | -0.20 | -0.10 |
| 0.00 | 18.18 | 18.04 | 18.15 | 18.60 |
| 0.10 | 16.19 | 15.98 | 16.06 | 16.42 |
| 0.20 | 14.80 | 14.82 | 14.74 | 15.08 |
| 0.30 | 13.79 | 13.75 | 13.90 | 14.32 |
| 0.40 | 13.30 | 13.24 | 13.44 | 13.60 |
| 0.50 | 13.18 | 13.09 | 13.28 | 13.48 |
| 0.60 | 13.24 | 13.17 | 13.29 | 13.49 |
| 0.70 | 13.48 | 13.52 | 13.61 | 13.72 |
| 0.80 | 13.98 | 14.03 | 14.05 | 14.25 |
| 0.90 | 14.89 | 14.89 | 14.95 | 15.05 |
| 1.00 | 16.40 | 16.31 | 16.32 | 16.46 |

**Fig. 3**. Dev set WERs for Density Ratio LM scaling factor $\lambda$ vs. sequence length scaling factor $\beta$. Here $\lambda = \lambda_\psi = \lambda_\tau$.

| WER | $\lambda\_\psi$ | | | |
|---|---|---|---|---|
| $\lambda\_\tau$ | 0.30 | 0.40 | 0.50 | 0.60 |
| 0.25 | 15.65 | | | |
| 0.30 | 14.32 | | | |
| 0.35 | 13.33 | 14.78 | | |
| 0.40 | 12.90 | 13.60 | | |
| 0.45 | 13.00 | 12.84 | 14.41 | |
| 0.50 | 13.58 | 12.63 | 13.48 | |
| 0.55 | 14.16 | 12.75 | 12.76 | 14.30 |
| 0.60 | | 13.20 | 12.66 | 13.49 |
| 0.65 | | 13.89 | 12.74 | 12.92 |
| 0.70 | | | 13.17 | 12.83 |
| 0.75 | | | 13.85 | 12.90 |
| 0.80 | | | | 13.32 |
| 0.85 | | | | 14.15 |

**Fig. 4**. Dev set WERs for different combinations of $\lambda_\tau$ and $\lambda_\psi$; sequence length scaling factor $\beta = -0.1$

.

of RNN-T, implemented as a non-blank symbol emission reward. (The blank symbol only consumes acoustic frames, not LM symbols [2]). Given that both factors have related effects on overall WER, the LM scaling factor(s) and the sequence length scaling factor need to be tuned jointly.

Fig. 2 and Fig. 3 illustrate the different relative sensitivities of WER to these factors for Shallow Fusion and the Density Ratio method, measured on the dev set.

In the second set of experiments, $\beta$ was fixed at -0.1, but the constraint $\lambda_\psi = \lambda_\tau$ was lifted, and a range of combinations was evaluated on the dev set. The results are shown in Fig. 4. The shading in Figs. 2, 3 and 4 uses the same midpoint value of 15.0 to highlight the results.

The best combinations of scaling factors from the dev set evaluations (see Fig. 2, Fig. 3 and Fig. 4) were used to generate the final eval set results, WERs and associated deletion, insertion and substitution rates, shown in Table 2. These results are summarized in Table 3, this time showing the exact values of LM scaling factor(s) used.

## 6. FINE-TUNING A YOUTUBE-TRAINED RNN-T USING LIMITED VOICE SEARCH AUDIO DATA

The experiments in Section 5 showed that an LM trained on text from the target Voice Search domain can boost the cross-domain performance of an RNN-T. The next experiments examined fine-tuning the original YouTube-trained RNN-T on varied, limited amounts of Voice Search {audio, transcript} data. After fine-tuning, LM fusion was applied, again comparing Shallow Fusion and the Density Ratio method.

Fine-tuning simply uses the YouTube-trained RNN-T model to warm-start training on the limited Voice Search

**Table 2**. In-domain and target domain performance of a YouTube-trained RNN-T, evaluated with and without fusion to a Voice Search LM (and normalizing YouTube LM in the case of the Density Ratio method).

| Model | WER | del | ins | sub |
|---|---|---|---|---|
| YouTube → YouTube | 11.3 | 2.4 | 1.9 | 7.0 |
| YouTube → Voice Search | 17.5 | 3.9 | 4.1 | 9.6 |
| Shallow Fusion | 14.5 | 4.6 | 4.1 | 5.8 |
| Density Ratio $\lambda_\psi = \lambda_\tau$ | 13.0 | 3.3 | 3.2 | 6.5 |
| Density Ratio $\lambda_\psi, \lambda_\tau$ | **12.5** | 3.9 | 2.9 | 5.7 |

{audio, transcript} data. This is an effective way of leveraging the limited Voice Search audio data: within a few thousand steps, the fine-tuned model reaches a decent level of performance on the fine-tuning task – though beyond that, it over-trains. A held-out set can be used to gauge over-training and stop training for varying amounts of fine-tuning data.

The experiments here fine-tuned the YouTube-trained RNN-T baseline using 10 hours, 100 hours and 1000 hours of Voice Search data, as described in Section 4.1. (The source domain RNN-LM was not fine-tuned). For each fine-tuned model, Shallow Fusion and the Density Ratio method were used to evaluate incorporation of the Voice Search RNN-LM, described in Section 5, trained on text transcripts from the much larger set of 21M Voice Search utterances. As in Section 5, the dev set was used to tune the LM scaling factor(s) and the sequence length scaling factor $\beta$. To ease parameter

5

tuning, the constraint $\lambda_\psi = \lambda_\tau$ was used for the Density Ratio method. The best combinations of scaling factors from the dev set were then used to generate the final eval results, which are shown in Table 3

**Table 3**. Fine tuning the YouTube-trained RNN-T baseline to the voice search target domain for different quantities of Voice Search fine-tuning data, evaluated with and without LM fusion on Voice Search test utterances. (Results for the "no fine-tuning" baseline carried over from Table 2).

| Model | WER | $\lambda$ | $\beta$ |
|---|---|---|---|
| Baseline (no fine-tuning) | 17.5 | - | -0.3 |
| Shallow Fusion | 14.5 | 0.3 | 0.6 |
| Density Ratio, $\lambda_\psi = \lambda_\tau$ | 13.0 | 0.5 | -0.3 |
| Density Ratio $\lambda_\psi, \lambda_\tau$ | **12.5** | 0.5, 0.6 | -0.1 |
| 10h fine-tuning | 12.5 | - | 0.0 |
| Shallow Fusion | 11.0 | 0.2 | 0.6 |
| Density Ratio, $\lambda_\psi = \lambda_\tau$ | 10.4 | 0.4 | 0.0 |
| Density Ratio, $\lambda_\psi, \lambda_\tau$ | **10.1** | 0.4, 0.45 | 0.0 |
| 100h fine-tuning | 10.6 | - | 0.0 |
| Shallow Fusion | 9.5 | 0.2 | 0.5 |
| Density Ratio, $\lambda_\psi = \lambda_\tau$ | **9.1** | 0.4 | 0.0 |
| 1,000h fine-tuning | 9.5 | - | 0.0 |
| Shallow Fusion | 8.8 | 0.2 | 0.5 |
| Density Ratio, $\lambda_\psi = \lambda_\tau$ | **8.5** | 0.3 | 0.0 |
| 21,000h fine-tuning | 7.8 | - | 0.1 |
| Shallow Fusion | 7.7 | 0.1 | 0.3 |
| Density Ratio, $\lambda_\psi = \lambda_\tau$ | **7.4** | 0.1 | 0.0 |

## 7. DISCUSSION

The experiments described here examined the generalization of a YouTube-trained end-to-end RNN-T model to Voice Search speech data, using varying quantities (from zero to 100%) of Voice Search audio data, and 100% of the available Voice Search text data. The results show that in spite of the vast range of acoustic and linguistic patterns covered by the YouTube-trained model, it is still possible to improve performance on Voice Search utterances significantly via Voice Search specific fine-tuning and LM fusion. In particular, LM fusion significantly boosts performance when only a limited quantity of Voice Search fine-tuning data is used.

The Density Ratio method consistently outperformed Shallow Fusion for the cross-domain scenarios examined, with and without fine-tuning to audio data from the target domain. Furthermore, the gains in WER over the baseline are significantly larger for the Density Ratio method than for Shallow Fusion, with up to 28% relative reduction in WER ($17.5 \rightarrow 12.5$) compared to up to 17% relative reduction ($17.5 \rightarrow 14.5$) for Shallow Fusion, in the no fine-tuning

scenario.

Notably, the "sweet spot" of effective combinations of LM scaling factor and sequence length scaling factor is significantly larger for the Density Ratio method than for Shallow Fusion (see Fig. 2 and Fig. 3). Compared to Shallow Fusion, larger absolute values of the scaling factor can be used.

A full sweep of the LM scaling factors ($\lambda_\psi$ and $\lambda_\tau$) can improve over the constrained setting $\lambda_\psi = \lambda_\tau$, though not by much. Fig. 4 shows that the optimal setting of the two factors follows a roughly linear pattern along an off-diagonal band.

Fine-tuning using transcribed Voice Search audio data leads to a large boost in performance over the YouTube-trained baseline. Nonetheless, both fusion methods give gains on top of fine-tuning, especially for the limited quantities of fine-tuning data. With 10 hours of fine-tuning, the Density Ratio method gives a 20% relative gain in WER, compared to 12% relative for Shallow Fusion. For 1000 hours of fine-tuning data, the Density Ratio method gives a 10.5% relative gave over the fine-tuned baseline, compared to 7% relative for Shallow Fusion. Even for 21,000 hours of fine-tuning data, i.e. the entire Voice Search training set, the Density Ratio method gives an added boost, from 7.8% to 7.4% WER, a 5% relative improvement.

A clear weakness of the proposed method is the apparent need for scaling factors on the LM outputs. In addition to the assumptions made (outlined in Section 3.1), it is possible that this is due to the implicit LM in the RNN-T being more limited than the RNN-LMs used.

## 8. SUMMARY

This article proposed and evaluated experimentally an alternative to Shallow Fusion for incorporation of an external LM into an end-to-end RNN-T model applied to a target domain different from the source domain it was trained on. The Density Ratio method is simple conceptually, easy to implement, and grounded in Bayes' rule, extending the classic hybrid ASR model to end-to-end models. In contrast, the most commonly reported approach to LM incorporation, Shallow Fusion, has no clear interpretation from probability theory. Evaluated on a YouTube $\rightarrow$ Voice Search cross-domain scenario, the method was found to be effective, with up to 28% relative gains in word error over the non-fused baseline, and consistently outperforming Shallow Fusion by a significant margin. The method continues to produce gains when fine-tuning to paired target domain data, though the gains diminish as more fine-tuning data is used. Evaluation using a variety of cross-domain evaluation scenarios is needed to establish the general effectiveness of the method.

### Acknowledgments

# 9. REFERENCES

[1] William Chan, Navdeep Jaitly, Quoc V. Le, and Oriol Vinyals, "Listen, attend and spell: A neural network for large vocabulary conversational speech recognition," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2016.

[2] Alex Graves, "Sequence transduction with recurrent neural networks," *CoRR*, vol. abs/1211.3711, 2012.

[3] Kartik Audhkhasi, Brian Kingsbury, Bhuvana Ramabhadran, George Saon, and Michael Picheny, "Building competitive direct acoustics-to-word models for english conversational speech recognition," in *Proc. IEEE ICASSP*, 2018.

[4] Chung-Cheng Chiu, Tara N. Sainath, Yonghui Wu, Rohit Prabhavalkar, Patrick Nguyen, Zhifeng Chen, Anjuli Kannan, Ron J. Weiss, Kanishka Rao, Ekaterina Gonina, Navdeep Jaitly, Bo Li, Jan Chorowski, and Michiel Bacchiani, "State-of-the-art speech recognition with sequence-to-sequence models," in *Proc. IEEE ICASSP*, 2018.

[5] Tomas Mikolov, Martin Karafiát, Lukás Burget, Jan Cernocký, and Sanjeev Khudanpur, "Recurrent neural network based language model," in *INTERSPEECH*, 2010.

[6] Jan Chorowski and Navdeep Jaitly, "Towards better decoding and language model integration in sequence to sequence models," in *Proc. Interspeech 2017*, 2017, pp. 523–527.

[7] Takaaki Hori, Shinji Watanabe, and John R. Hershey, "Multi-level language modeling and decoding for open vocabulary end-to-end speech recognition," in *Proc. ASRU*, 2017.

[8] Anjuli Kannan, Yonghui Wu, Patrick Nguyen, Tara Sainath, Zhifeng Chen, and Rohit Prabhavalkar, "An analysis of incorporating an external language model into a sequence-to-sequence model," in *Proc. IEEE ICASSP*, 2018.

[9] Shubham Toshniwal, Anjuli Kannan, Chung-Chen Chiu, Yonghui Wu, Tara Sainath, and Karen Livescu, "A comparison of techniques for language model integration in encoder-decoder speech recognition," *CoRR*, 2018.

[10] Caglar Gulcehre, Orhan Firat, Kelvin Xu, Kyunghyun Cho, Loic Barrault, Huei-Chi Lin, Fethi Bougares, Holger Schwenk, and Yoshua Bengio, "On using monolingual corpora in neural machine translation," *CoRR*, vol. abs/1503.03535, 2015.

[11] Anuroop Sriram, Heewoo Jun, Sanjeev Satheesh, and Adam Coates, "Cold fusion: Training seq2seq models together with language models," in *Proc. Interspeech*. 2018, ISCA.

[12] Changhao Shan, Chao Weng, Guangsen Wang, Dan Su, Min Luo, Dong Yu, and Lei Xie, "Component fusion: learning replaceable language model component for end-to-end speech recognition system," in *Proc. IEEE ICASSP*, 2019.

[13] Frederick Jelinek, "Continuous speech recognition by statistical methods," in *Proceedings of the IEEE 64*, 1976.

[14] Yasuhiro Minami, Toshiyuki Hanazawa, Hitoshi Iwamida, Erik McDermott, Kiyohiro Shikano, Shigeru Katagiri, and Masaona Nakagawa, "On the Robustness of HMM and ANN Speech Recognition Algorithms," in *Proc. International Conference on Spoken Language Processing*, 1990, vol. 2, pp. 1345–1348.

[15] Patrick Haffner, "Connectionist speech recognition with a global MMI algorithm," in *Proc. Eurospeech*. 1993, ISCA.

[16] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.

[17] Brian Kingsbury, "Lattice-based optimization of sequence classification criteria for neural-network acoustic modeling," in *ICASSP*, 2009, pp. 3761–3764.

[18] Frank Seide, Gang Li, Xie Chen, and Dong Yu, "Feature engineering in context-dependent deep neural networks for conversational speech transcription," in *Proc. ASRU*, 2011.

[19] Sepp Hochreiter and Juergen Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, Nov. 1997.

[20] Hasim Sak, Oriol Vinyals, Georg Heigold, Andrew W. Senior, Erik McDermott, Rajat Monga, and Mark Z. Mao, "Sequence discriminative distributed training of long short-term memory recurrent neural networks," in *Proc. Interspeech*. 2014, ISCA.

[21] Hervé Bourlard and Nelson Morgan, *Connectionist speech recognition - A Hybrid Approach*, Kluwer Academic Publishers, 1994.

[22] Ehsan Variani, Tom Bagby, Erik McDermott, and Michiel Bacchiani, "End-to-end training of acoustic models for large vocabulary continuous speech recognition with tensorflow," in *Proc. Interspeech*. 2018, ISCA.

[23] Mirco Ravanelli, Titouan Parcollet, and Yoshua Bengio, "The pytorch-kaldi speech recognition toolkit," in *ICASSP*, 2019, pp. 6465–6469.

[24] Alex Graves, Santiago Fernndez, and Faustino Gomez, "Connectionist temporal classification: Labelling unsegmented sequence data with recurrent neural networks," in *In Proceedings of the International Conference on Machine Learning, ICML 2006*, 2006, pp. 369–376.

[25] Takaaki Hori, Jaejin Cho, and Shinji Watanabe, "End-to-end speech recognition with word-based RNN language models," *CoRR*, vol. abs/1808.02608, 2018.

[26] Naoyuki Kanda, Xugang Lu, and Hisashi Kawai, "Maximum a posteriori based decoding for end-to-end acoustic models," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 25, no. 05, pp. 1023–1034, Aug. 2017.

[27] Erik McDermott, "A deep generative acoustic model for compositional automatic speech recognition," in *Proceedings of Neural Information Processing Systems (NeurIPS) Workshop: Interpretability and Robustness in Audio, Speech, and Language*, 2018.

[28] L. R. Rabiner and B.-H. Juang, *Fundamentals of Speech Recognition*, PTR Prentice-Hall, Inc., Englewood Cliffs, New Jersey 07632, 1993.

[29] Hank Liao, Erik McDermott, and Andrew W. Senior, "Large scale deep neural network acoustic modeling with semi-supervised training data for youtube video transcription," *2013 IEEE Workshop on Automatic Speech Recognition and Understanding*, pp. 368–373, 2013.

[30] Hagen Soltau, Hank Liao, and Hasim Sak, "Neural speech recognizer: Acoustic-to-word LSTM model for large vocabulary speech recognition," *CoRR*, vol. abs/1610.09975, 2016.

[31] Sami Virpioja, Peter Smit, Stig-Arne Grnroos, and Mikko Kurimo, "Morfessor 2.0: Python implementation and extensions for morfessor baseline," D4 julkaistu kehittmis- tai tutkimusraportti tai -selvitys, 2013.

[32] Diederik Kingma and Jimmy Ba, "Adam: A method for stochastic optimization," in *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.