# A Distributed Optimisation Framework Combining Natural Gradient with Hessian-Free for Discriminative Sequence Training

Adnan Haider[1,2], Chao Zhang[2], Florian L. Kreyssig[3] and Philip C. Woodland*

*Cambridge University Engineering Department, Trumpington Street, Cambridge, CB2 1PZ U.K.*

ABSTRACT

This paper presents a novel natural gradient and Hessian-free (NGHF) optimisation framework for neural network training that can operate efficently in a distributed manner. It relies on the linear conjugate gradient (CG) algorithm to combine the natural gradient (NG) method with local curvature information from Hessian-free (HF) or other second-order methods. A solution to a numerical issue in CG allows effective parameter updates to be generated with far fewer CG iterations than usually used (*e.g.* 5-8 instead of 200). This work also presents a novel preconditioning approach to improve the progress made by individual CG iterations for models with shared parameters. Although applicable to other training losses and model structures, NGHF is investigated in this paper for lattice-based discriminative sequence training for hybrid hidden Markov model acoustic models using a standard recurrent neural network, long short-term memory, and time delay neural network models for output probability calculation. Automatic speech recognition experiments are reported on the multi-genre broadcast data set for a range of different acoustic model types. These experiments show that NGHF achieves larger word error rate reductions than standard stochastic gradient descent or Adam, while requiring orders of magnitude fewer parameter updates.

## 1. Introduction

With the availability of increased computing power and appropriate parameter initialisation methods (Hinton et al., 2006, 2012; Glorot and Bengio, 2010), the hybrid automatic speech recognition (ASR) approach, *i.e.* the use of hidden Markov model (HMM) acoustic model with neural networks (NNs) for observation probability calculation, has achieved state-of-the-art performance on large vocabulary continuous speech recognition (LVCSR) tasks (Xiong et al., 2016; Saon et al., 2016). Whilst NNs can model the underlying non-linear manifold of speech data well (Hinton et al., 2012), their deep structures create complex dependencies among parameters that can make these models difficult to train with standard, or variants of, stochastic gradient descent (SGD) (Desjardins et al., 2015; Grosse and Salakhudinov, 2015).

Due to the sequential nature of speech, sequence-level training objectives are often used for ASR. It has been observed that training acoustic models using a discriminative loss function that not only maximises the probability of the reference transcription or recognition accuracy but also minimises that of all competing hypotheses often improves ASR performance by a significant margin (Valtchev et al., 1997; Povey, 2005; Graves and Jaitly, 2014; Povey et al., 2016). Therefore discriminative loss functions have become a standard choice for state-of-the-art ASR systems (Saon et al., 2016; Xiong et al., 2016; Chiu et al., 2018; Lüscher et al., 2019). For training NNs with such a loss, SGD with training data shuffled at the utterance level is widely used. Often

this use of SGD requires considerable skill in choosing suitable values for hyper-parameters such as the learning rate or the momentum coefficient, while second-order optimisation methods require fewer hyper-parameters and can lead to more effective parameter updates compared to SGD (Xu et al., 2020). These approaches leverage the local curvature information contained in the Hessian matrix (*i.e.* the second-order derivatives w.r.t. the model parameters) to overcome the issues that SGD encounters with highly non-linear and ill-conditioned loss functions (Becker and LeCun, 1988). However, a major weakness of second-order methods is that they are computationally inefficient for models with a large number of parameters. This is because scaling the gradient directions by the inverse of the Hessian matrix requires computation on the order of $\mathbb{O}(D^3)$ where $D$ is the number of model parameters.

A large batch second-order approach, the Hessian-free (HF) method, was proposed by Martens (2010). Instead of computing the Hessian matrix explicitly, HF finds an approximate solution by solving an equivalent system using the conjugate gradient (CG) algorithm. Compared to SGD with mini-batches typically containing only a few utterances, HF can yield more effective updates based on very large batches (*e.g.* thousands of utterances ) and requires far fewer parameter updates to converge. Hence, it is easier to parallelise HF than SGD by distributing the training data in each large batch among all available computation units and accumulating the outcomes to yield a parameter update. However, in practice, HF can require more computation than SGD, since it often needs a large number (*e.g.* 200) of CG iterations for each parameter update (Martens, 2010; Kingsbury et al., 2012; Wiesler et al., 2013; Martens and Grosse, 2015).

Recently, the natural gradient (NG) method, which was originally proposed for training small models with maximum likelihood (ML) (Amari, 1997), has received renewed

---

*Corresponding author.

✉ adnan_haider@apple.com (A. Haider); cz277@cam.ac.uk (C. Zhang); flk24@cam.ac.uk (F.L. Kreyssig); pcw@eng.cam.ac.uk (P.C. Woodland)

[1]Work was done while Adnan Haider was at Cambridge University
[2]Equal contributions.
[3]Funded by an EPSRC Doctoral Training Partnership Award.

popularity for training deep NN models (Pascanu and Bengio, 2013; Desjardins et al., 2015). NG computes the steepest descent direction in the space of the model output distributions instead of the space of the model parameters, and can result in fewer generalisation errors and faster convergence (Roux et al., 2008; Bernacchia et al., 2018). NG requires calculating the product of the inverse of the Fisher information (FI) matrix with the gradient estimate. As for second-order approaches, directly calculating the inverse of the FI matrix has a $\mathbb{O}(D^3)$ computational complexity, rendering it unusable for large models. To reduce this cost, approaches such as enforcing a block diagonal structure on the FI matrix and approximating the diagonal blocks as the Kronecker product of two smaller matrices have been introduced (Povey et al., 2015; Martens and Grosse, 2015; Grosse and Salakhudinov, 2015; George et al., 2018).

In this paper, we propose an approach to stabilise the computation of the directional derivative that alleviates the numerical instability that can accompany CG when applied to NNs. To improve training on models with shared architectures, this work proposes a novel a preconditioning approach to speed up the progress made by CG. Furthermore, in order to estimate the inverse of the FI matrix without making any assumptions about its structure, we propose applying CG to NG for discriminative sequence training. As a result, CG can provide a common optimisation framework to use NG jointly with HF or other second-order approaches. We, therefore, propose the NGHF method that combines the advantages of NG and HF and further improves the stability and performance of NN model training. The efficacy of the proposed optimisation framework and the NGHF method is evaluated by conducting ASR experiments on a 200 hour multi-genre broadcast (MGB) speech recognition dataset (Bell et al., 2015; Woodland et al., 2015). Discriminative sequence training with hybrid HMM acoustic models is studied, which use long short-term memory networks (LSTMs) (Hochreiter and Schmidhuber, 1997), recurrent neural networks (RNNs) (Rumelhart et al., 1986; Elman, 1990), or time-delayed neural networks (TDNNs) (Waibel et al., 1989; Peddinti et al., 2015) for output probability calculations. Compared to SGD and Adam (Kingma and Ba, 2015), NGHF is demonstrated to have a faster convergence speed and can result in lower word error rates (WERs).

This paper is organised as follows. Section 2 and 3 review related work on training NNs and the HF method with CG. Section 4 presents our improved implementation of CG along with the distributed optimisation framework. Sections 5 and 6 describe the proposed CG-based NG and NGHF methods for discriminative sequence training. The experimental setup and results are given in Sec. 7 and Sec. 8, followed by a discussion in Sec. 9 and finally conclusions.

Compared to our previous papers (Haider and Woodland, 2017, 2018), this paper proposes a more general approach to regulate the scaling and rotation induced on the gradient direction during NG descent learning with second-order methods (Sec. 6). All techniques are described more completely and comprehensively. The design of the distributed CG-based optimisation method is fully presented (Sec. 4.1). The description and analysis of how to improve the numerical stability of CG (Sec. 4.2) and a preconditioning method that improves the performance for models with shared parameters (Sec. 4.3) are included. The experimental results have been extended and now include comparisons of the proposed method for TDNNs, LSTMs and other recurrent models.

## 2. Related work

Starting from layer-by-layer pre-training and improved random initialisation methods (Hinton and Salakhutdinov, 2006; Glorot and Bengio, 2010), optimisation improvements have played a critical role in the development of deep learning. Various methods have been applied to improve the stability and efficiency for training complex NN structures with SGD. Even though the use of these methods has significantly improved NN training, the task of finding a more suitable optimisation procedure for supervised NN training continues to be an active area of research. This section first reviews different SGD, HF, and NG methods, with a focus on distributed computation. Then discriminative sequence training is reviewed, which is a key focus of this paper.

### 2.1. Synchronous and Asynchronous SGD

Adapting the learning rates during training is known to improve the stability and convergence speed of SGD. Rule-based strategies can be used to decay the learning rate shared by all DNN parameters (Renals et al., 1992; Bottou, 2010; Senior et al., 2013). Methods such as RMSProp (Tieleman and Hinton, 2012), AdaGrad (Duchi et al., 2011), AdaDelta (Zeiler, 2012), and Adam (Kingma and Ba, 2015) assign a separate learning rate to each model parameter. These are updated based on previous gradients. Parameter updates can also be derived as the output of an RNN based on the gradients (Andrychowicz et al., 2016). The methods developed in this paper do not require tuning of learning rates and implicitly assign separate learning rates to each parameter.

In standard SGD, greater parallelisation is achieved by using a larger mini-batch. The (equivalent) mini-batch size can be increased in a synchronous manner, which requires every "worker" process to calculate the gradients for their mini-batch for the same copy of the model parameters. These are then combined (Zinkevich et al., 2010). The main issue of synchronous distributed SGD is the high transmission cost (Seide et al., 2014; Ström, 2015). Asynchronous SGD (ASGD) can reduce the transmission cost by allowing each worker to keep a distinct copy of the model and to communicate independently with a centralised parameter server (Dean et al., 2012; Heigold et al., 2014) or other workers (Zhang et al., 2019b,c). This also increases the number of updates. The blockwise model update filtering method divides workers into blocks and averages the models produced by each block with a mechanism analogous to the use of momentum in standard SGD (Chen and Huo, 2016; Ladkat et al., 2019). The optimisation methods developed in this paper can be used in a highly parallel and distributed man-

ner, which relies on the use of very large batches and a small number of updates to efficiently and accurately exploit parallelisation.

## 2.2. HF and NG based Approaches

The HF optimisation framework was initially proposed to minimise the squared errors for auto-encoders based on the linear CG algorithm (Shewchuk, 1994) and the Gauss-Newton (GN) approximation (Martens, 2010). Later studies applied the approach to cross-entropy (CE) training for small scale ASR and hand-written digit classification (Vinyals and Povey, 2012; Wiesler et al., 2013). Kingsbury et al. (2012) extended HF to lattice-based discriminative sequence training for LVCSR with a large batch size, which allows the data to be processed in parallel based on a master/worker structure. The workers compute the gradients and the curvature information based on the same copy of the model parameters in a distributed fashion, while the master process collects the outputs from all workers to perform the CG algorithm to update the model parameters. Momentum and preconditioning can be integrated into this optimisation framework as shown in later studies (Sainath et al., 2013a,b). To reduce the computation workload of the master, only a small percentage (*e.g.* 1%) of the training data is sampled for use in each CG iteration (Martens, 2010; Kingsbury et al., 2012).

Compared to the original work on NG that is aimed at training small models with the ML criterion (Amari, 1997), more recent work focuses on training large NN models with tens of millions of parameters, in which case it is computationally impractical to estimate the inverse of the FI matrix. To overcome this issue, most studies assume that the FI matrix has some form of block diagonal structure (Roux et al., 2008; Povey et al., 2015; Martens and Grosse, 2015; Grosse and Salakhudinov, 2015; George et al., 2018). Kronecker-factored approximate curvature (K-FAC) is such an approach that assumes the parameters from different layers are independent and approximates each block of the FI matrix as the Kronecker product of two smaller matrices (Martens and Grosse, 2015). Povey et al. (2015) proposed a similar idea and applied it to lattice-free maximum mutual information training (Povey et al., 2016), which was observed to yield no performance loss when simply averaging the models produced by the different workers configured to operate asynchronously. Interestingly, the second-order momentum used in Adam can be viewed as an approximate diagonal FI matrix that assumes every model parameter is independently estimated (Kingma and Ba, 2015). Schulman et al. (2015) proposed using CG to estimate the inverse of the FI matrix without enforcing any assumptions for reinforcement learning. Haider and Woodland (2017) proposed a similar method for discriminative sequence training for ASR.

## 2.3. Discriminative Sequence Training

An ASR system aims to convert the acoustic feature sequence $\mathbf{O}$ of a speech utterance to its underlying word sequence $\mathbf{W}^{\text{ref}}$. In practice this is often achieved by finding the hypothesised word sequence $\hat{\mathbf{W}}$ that gives the maximum posterior probability $P(\mathbf{W}|\mathbf{O})$, *i.e.*

$$\hat{\mathbf{W}} = \arg\max_{\mathbf{W}} P(\mathbf{W}|\mathbf{O}). \tag{1}$$

In the noisy source-channel framework for ASR, $P(\mathbf{W}|\mathbf{O})$ is calculated using an HMM-based acoustic model that produces $p(\mathbf{O}|\mathbf{W}, \theta)$ (denoted as $p_\theta(\mathbf{O}|\mathbf{W})$ in this paper) and a language model (LM) that estimates $P(\mathbf{W})$. In practice, the LM is first trained separately and then combined at test-time to decode the input acoustic feature vectors for an utterance into the most probable word sequence.

Let $\mathcal{M}$ denote the space of all probability distributions $P_\theta(\mathbf{W}|\mathbf{O})$ that can result from different parameter configurations of a chosen NN when employed within an HMM. The goal of learning is to identify a viable candidate in $\mathcal{M}$ that achieves the greatest reduction in the empirical loss, which refers to the average loss over all training samples w.r.t. a given risk function while generalising well to new examples. To effectively train acoustic models, two forms of discriminative sequence level criterion are commonly used. The first form corresponds to the maximum mutual information (MMI) loss. For a single observation feature sequence $\mathbf{O}$, the MMI loss corresponds to

$$\mathcal{L}_{\text{MMI}}(\theta) = -\ln \frac{p_\theta(\mathbf{O}|\mathbf{W}^{\text{ref}})^\kappa P(\mathbf{W}^{\text{ref}})}{\sum_{\mathbf{W}} p_\theta(\mathbf{O}|\mathbf{W})^\kappa P(\mathbf{W})}, \tag{2}$$

The method was initially proposed for spoken digit recognition (Bahl et al., 1986), where $\kappa$ is the acoustic scaling factor that re-scales the acoustic model and language model scores to be in the same range (Woodland and Povey, 2002). The MMI loss can be viewed as maximising the probability of $\mathbf{W}^{\text{ref}}$ while also minimising that of every competing hypothesis $\mathbf{W}$, and is thus a discriminative sequence training loss. When applied to LVCSR (Valtchev et al., 1997), the denominator of Eqn. (2) needs to be calculated efficiently and often relies on using word lattices for each utterance as a compact representation of all of the important competing hypotheses.

The minimum Bayes risk (MBR) loss is another commonly used loss function for discriminative sequence training (Goel and Byrne, 2000; Kaiser et al., 2000), which directly minimises ASR errors by using WER related metrics as the risk function. The MBR loss for a single observation feature sequence $\mathbf{O}$ is defined as

$$\mathcal{L}_{\text{MBR}}(\theta) = \frac{\sum_{\mathbf{W}} p_\theta(\mathbf{O}|\mathbf{W})^\kappa P(\mathbf{W})\mathcal{A}(\mathbf{W}, \mathbf{W}^{\text{ref}})}{\sum_{\mathbf{W}} p_\theta(\mathbf{O}|\mathbf{W})^\kappa P(\mathbf{W})}, \tag{3}$$

where $\mathcal{A}(\mathbf{W}, \mathbf{W}^{\text{ref}})$ is the risk function measuring the difference between the reference and a competing hypothesis. For ASR, approximations to the number of word-level, phone-level, and HMM state-level errors are widely used as the risk function (Povey and Woodland, 2002; Gibson and Hain, 2006; Shannon, 2017). When phone-level errors are used, the MBR loss is called minimum phone error (MPE) (Povey and Woodland, 2002), which will be the MBR loss used in the experiments in this paper. As for MMI, for LVCSR word lattices of each training utterance are required since the calculation of the MBR loss involves all competing hypotheses.

Lattice-based discriminative sequence training with an MMI or MBR loss has also been widely used to finetune NN-HMM hybrid systems (Valtchev, 1995; Kingsbury, 2009; Veselỳ et al., 2013; Su et al., 2013; Wiesler et al., 2015; Zhang and Woodland, 2015), which are often initialised by frame-level training with the CE loss (Hinton et al., 2012). Lattice-free MMI uses a general phone-level recognition network to replace utterance-specific lattices, which enables efficient processing of multiple training utterances in parallel (Povey et al., 2016).

## 3. Hessian-free Optimisation

This section provides an overview of the HF optimisation framework (Martens, 2010; Kingsbury et al., 2012). At the core of all first and second-order optimisation methods is Taylor's theorem. Assuming the loss function $\mathcal{L}(\theta)$ is sufficiently smooth, the *second-order Taylor approximation* employs the following quadratic function to locally approximate the function as

$$\mathcal{L}(\theta + \Delta\theta) \approx \mathcal{L}(\theta) + \Delta\theta^{\mathrm{T}}\nabla_\theta\mathcal{L}(\theta) + \frac{1}{2}\Delta\theta^{\mathrm{T}}\mathbf{H}\Delta\theta, \quad (4)$$

where $\nabla_\theta$ is the gradient operator in the space of $\theta$, $\Delta\theta$ represents an offset within a convex neighbourhood of $\theta$, and $\mathbf{H}$ is the Hessian matrix of $\mathcal{L}$ w.r.t. $\theta$, *i.e.* $\mathbf{H} = \nabla_\theta^2\mathcal{L}(\theta)$.

Instead of optimising the loss function directly, at each iteration of the optimisation process, second-order methods focus on a generating a candidate update $\Delta\theta$ through minimising Eqn. (4) where $\mathbf{H}$ is approximated by a candidate matrix $\mathbf{B}$. Differentiating Eqn. (4) and setting it to zero yields the *Newton direction*

$$\Delta\theta = -\mathbf{B}^{-1}\nabla_\theta\mathcal{L}(\theta). \quad (5)$$

However, computing this direction directly is expensive since it requires $\mathbb{O}(D^2)$ complexity to store $\mathbf{B}$ and $\mathbb{O}(D^3)$ to invert it. These obstacles, however, can be overcome by employing inexact Newton methods such as the CG algorithm.

---

**Algorithm 1** The linear conjugate gradient (CG) algorithm.

Let $M$ be the number of CG iterations to execute
Set $\boldsymbol{v}_0 \leftarrow -\nabla_\theta\mathcal{L}(\theta), \boldsymbol{r}_0 \leftarrow \boldsymbol{v}_0, m \leftarrow 0$
**while** $m < M$ **do**
    Compute $\boldsymbol{r}_m^{\mathrm{T}}\boldsymbol{r}_m$
    Set $\alpha_m \leftarrow \boldsymbol{r}_m^{\mathrm{T}}\boldsymbol{r}_m / \boldsymbol{v}_m^{\mathrm{T}}\mathbf{B}\boldsymbol{v}_m$
    Update $\Delta\theta_{m+1} \leftarrow \Delta\theta_m + \alpha_m\boldsymbol{v}_m$
    Update $\boldsymbol{r}_{m+1} \leftarrow \boldsymbol{r}_m - \alpha_m\mathbf{B}\boldsymbol{v}_m$
    Compute $\boldsymbol{r}_{m+1}^{\mathrm{T}}\boldsymbol{r}_{m+1}$
    Set $\beta_{m+1} \leftarrow \boldsymbol{r}_{m+1}^{\mathrm{T}}\boldsymbol{r}_{m+1} / \boldsymbol{r}_m^{\mathrm{T}}\boldsymbol{r}_m$
    Update $\boldsymbol{v}_{m+1} \leftarrow \boldsymbol{r}_{m+1} + \beta_{m+1}\boldsymbol{v}_m$
**return** $\Delta\theta$ as the one that leads to the best performance on the validation set among $\Delta\theta_1, \Delta\theta_2, \dots, \Delta\theta_M$

---

### 3.1. The CG Algorithm

CG is an iterative algorithm that implicitly minimises the quadratic function

$$g(\Delta\theta_m) = \frac{1}{2}\Delta\theta_m^{\mathrm{T}}\mathbf{B}\Delta\theta_m + \Delta\theta_m^{\mathrm{T}}\nabla_\theta\mathcal{L}(\theta) \quad (6)$$

by solving the linear linear system

$$\mathbf{B}\Delta\theta = -\nabla_\theta\mathcal{L}(\theta), \quad (7)$$

At each iteration $m$, the algorithm minimises Eqn. (6) by taking an appropriate step size $\alpha_m$ along a conjugate search direction $\boldsymbol{v}_m$ w.r.t. $\mathbf{B}$ such that the direction is never revisited at subsequent iterations. When $\mathbf{B}$ is symmetric and positive definite, the solution to the linear system yields a unique minimiser of Eqn. (6). Since Eqn. (6) only approximates $\mathcal{L}(\theta)$, the standard practice in training parametric models such as NN is to run only finite iterations of the algorithm in which minimising the quadratic function correlates with reductions in the empirical loss (Martens, 2020).

The detailed CG procedure is presented as Algorithm 1, for which an excellent explanation can be found in (Shewchuk, 1994). The key features of CG are summarised below.

- $\boldsymbol{v}_1, \boldsymbol{v}_2, \dots, \boldsymbol{v}_M$ are $\mathbf{B}$-orthogonal. This means any direction $\boldsymbol{v}_m$ is conjugate to any other direction w.r.t. $\mathbf{B}$.

- CG computes $\mathbf{B}\boldsymbol{v}_m$ instead of the Hessian matrix itself. When $\mathbf{B}$ is chosen to approximate the Hessian matrix, the method is known as Hessian free (Martens, 2010; Bottou et al., 2018).

- Since $\Delta\theta_M = -\alpha_0\nabla_\theta\mathcal{L}(\theta) + \sum_{m=1}^{M-1}\alpha_m\boldsymbol{v}_m$, $\Delta\theta_0$ equals $-\alpha_0\nabla_\theta\mathcal{L}(\theta)$ and can be seen as the update obtained by gradient descent with an optimal learning rate that minimises Eqn. (6).

- CG will converge monotonically to the exact Newton direction within $M$ iterations, if $\mathbf{B}$ has $M$ distinct or clustered eigenvalues (Nocedal and Wright, 2016).

### 3.2. Approximating the Hessian with the Gauss-Newton Matrix

This section reviews using the GN matrix as $\mathbf{B}$ in the approximation of $\mathbf{H}$ in the HF method (Martens and Sutskever, 2011). For simplicity, here we consider the case with only one input sample, a frame at time $t$. It is straightforward to generalise the method and equations to the case of many input samples. Let $\boldsymbol{a}_t^{\mathrm{out}} = \{a_{t,1}, a_{t,2}, \dots, a_{t,K}\}$ be the logit values (the input values to the softmax output activation function), $K$ be the output layer size, $\mathbf{H}_{ij}$ be the element of the $i$ th row and $j$ th column of $\mathbf{H}$, $\mathbf{H}_{ij}$ can be written as

$$\mathbf{H}_{ij} = \frac{\partial}{\partial\theta_j}\left(\frac{\partial\mathcal{L}(\theta)}{\partial\theta_i}\right) \quad (8)$$

$$= \sum_{k=1}^{K}\frac{\partial a_{t,k'}}{\partial\theta_j}\sum_{k'=1}^{K}\frac{\partial a_{t,k}}{\partial\theta_i}\frac{\partial^2\mathcal{L}(\theta)}{\partial a_{t,k}\partial a_{t,k'}} + \sum_{k=1}^{K}\frac{\partial\mathcal{L}(\theta)}{\partial a_{t,k}}\frac{\partial^2 a_{t,k}}{\partial\theta_i\partial\theta_j}$$

In Eqn. (8), the first term can be interpreted as the contribution to the Hessian made by the variation in $a_t^{\text{out}}$, while the second term is the contribution due to the variation in $\boldsymbol{\theta}$. If $\boldsymbol{\theta}$ is around a region of local minimum w.r.t. the average empirical loss over samples, then $\partial \mathcal{L}(\boldsymbol{\theta})/\partial a_{t,k} \approx 0$ and the second term is negligible. As a result,

$$\mathbf{H}_{ij} \approx \sum_{k=1}^{K} \frac{\partial a_{t,k'}}{\partial \theta_j} \sum_{k'=1}^{K} \frac{\partial a_{t,k}}{\partial \theta_i} \frac{\partial^2 \mathcal{L}(\boldsymbol{\theta})}{\partial a_{t,k} \partial a_{t,k'}} = \mathbf{G}_{ij}, \qquad (9)$$

which is an element of the GN matrix, $\mathbf{G}$. By rearranging Eqn. (9), the GN matrix can be written as

$$\mathbf{G} = \mathbf{J}^{\text{T}} (\nabla^2_{a_t^{\text{out}}} \mathcal{L}(\boldsymbol{\theta})) \, \mathbf{J}, \qquad (10)$$

where $\mathbf{J}$ is the Jacobian matrix $\nabla_{\boldsymbol{\theta}}(a_t^{\text{out}})$, and $\nabla^2_{a_t^{\text{out}}} \mathcal{L}(\boldsymbol{\theta})$ is the Hessian matrix w.r.t. $a_t^{\text{out}}$.

As shown by Schraudolph (2002), for matching loss functions where $\nabla^2_{a_t^{\text{out}}} \mathcal{L}(\boldsymbol{\theta})$ is positive definite w.r.t. the NN output logits, the GN matrix is guaranteed to be positive semi-definite. When the GN approximation is applied to to lattice-based MBR training (Kingsbury et al., 2012), the component $\nabla^2_{a_t^{\text{out}}} \mathcal{L}_{\text{MBR}}(\boldsymbol{\theta})$ which we denote as $\hat{\mathbf{H}}$ takes the following form (Haider, 2019):

$$\hat{\mathbf{H}} = \kappa^2 \left( \text{diag}(\boldsymbol{\gamma}_t^{\text{MBR}}) - \boldsymbol{\gamma}_t^{\text{MBR}} \boldsymbol{\gamma}_t^{\text{MBR T}} \right), \qquad (11)$$

where $\kappa$ is the acoustic scaling factor, $\gamma_{t,k}$ and $\gamma_{t,k}^{\text{MBR}}$ are correspondingly the ML and MBR occupancy at time $t$ w.r.t. the HMM state $k$ (tied to the DNN output unit $k$), and $\text{diag}(\cdot)$ converts a vector into a diagonal matrix.

Regarding ML training in the ASR literature, the term "occupancy" often refers to $\gamma_{t,k} = P_{\boldsymbol{\theta}}(q_t = k|\mathbf{O}, \mathbf{W}^{\text{ref}})$, the posterior probability showing how probable frame $t$ is aligned with state $k$ given a pair of sequences $\mathbf{O}$ and $\mathbf{W}^{\text{ref}}$. $\gamma_{t,k}$ is often calculated using the *forward-backward* procedure (Baum and Eagon, 1967). Regarding MBR sequence training, the occupancy $\gamma_{t,k}^{\text{MBR}}$ is defined based on the gradients of the loss function $\partial \mathcal{L}_{\text{MBR}}(\boldsymbol{\theta})/\partial a_{t,k} = -\kappa \gamma_{t,k}^{\text{MBR}}$, where $\gamma_{t,k}^{\text{MBR}} = \gamma_q (c_q - c_{\text{avg}}) \gamma_{t,k}$, $\gamma_q$ is the occupancy passing through arc $q$, $c_{\text{avg}}$ and $c_q$ are the weighted average correctness of all hypotheses and the hypotheses including arc $q$. $\gamma_q$, $c_q$, and $c_{\text{avg}}$ can be collected by performing a modified forward-backward procedure to align every arc in the lattice with $\mathbf{O}$ (Povey, 2005).

From Eqn. (11), it can seen that $\nabla^2_{a_t^{\text{out}}} \mathcal{L}_{\text{MBR}}(\boldsymbol{\theta})$ is not positive definite w.r.t. the NN output logits and hence the GN matrix is no longer guaranteed to be positive semi-definite. Interestingly, even through the GN matrix no longer possesses the property of being positive semi-definite, its use as an approximation to the Hessian has been shown empirically to be effective in obtaining stable WER reductions from lattice-based MBR sequence training (Kingsbury et al., 2012; Sainath et al., 2013a,b; Dognin and Goel, 2013). The next section describes a particular property of the GN matrix that has been recently shown by Haider and Woodland

(2018) to be effective when using a disciminative sequence criterion to train NN models with sharp softmax distributions.

### 3.3. Scaling by the Gauss-Newton Matrix

The frame-level CE loss used to train an NN acoustic model from random initialisation using frame-to-HMM-state alignments often results in distribution of the NN softmax outputs being very sharp, in particular for models that use ReLU activation functions. This was reported in (Haider, 2019) to be a contributing factor in achieving only very small improvements from discriminative sequence training. A similar issue has been observed in connectionist temporal classification (Graves et al., 2006), where the sharp distributions are caused by the blank unit instead of the 0-1 training labels. According to Eqn. (9), $\mathbf{G}$ captures the curvature of the training loss w.r.t. the model output distribution. Re-scaling the gradients by $\mathbf{G}^{-1}$ effectively de-weights the back-propagated information that can induce large changes in the loss value in EBP. In the context of discriminative sequence training, this regularises the sharp model output distributions and improves the performance of MBR training (Haider, 2019).

### 3.4. Matrix-Vector-Products with the Gauss-Newton Matrix

Within each iteration of CG, a multiplication of the GN matrix with a vector $\boldsymbol{v}$ ($\mathbf{B} \boldsymbol{v}_m$ in Alg. 1) in the parameter space, $\mathbf{G} \boldsymbol{v} = \mathbf{J}^{\text{T}} (\nabla^2_{a_t^{\text{out}}} \mathcal{L}(\boldsymbol{\theta})) \mathbf{J} \boldsymbol{v}$, corresponds to the following sequential multiplication:

- Computing the directional derivative $\mathbf{J} \boldsymbol{v}$ using a modified forward propagation procedure.

- Multiplying the resulting vector $\mathbf{J} \boldsymbol{v}$ by $\nabla^2_{a_t^{\text{out}}} \mathcal{L}(\boldsymbol{\theta})$ which corresponds to $\hat{\mathbf{H}}$ in this context.

- Since the error backpropagation (EBP) procedure computes $\mathbf{J}^{\text{T}} (\nabla_{a_t^{\text{out}}} \mathcal{L}(\boldsymbol{\theta}))$, $\mathbf{G} \boldsymbol{v}$ can be obtained using EBP by replacing $\nabla_{a_t^{\text{out}}} \mathcal{L}(\boldsymbol{\theta})$ with $(\nabla^2_{a_t^{\text{out}}} \mathcal{L}(\boldsymbol{\theta})) \mathbf{J} \boldsymbol{v}$.

To compute $\mathbf{J} \boldsymbol{v}$ efficiently using a modified forward propagation procedure, Pearlmutter (1994) introduced an operator $\mathcal{R}(\cdot)$ to calculate the directional derivative $\nabla_{\boldsymbol{\theta}}(\cdot) \boldsymbol{v}$, and there is $\mathcal{R}(\boldsymbol{\theta}) = \boldsymbol{v}$. For a fully-connected (FC) layer with $a_{t,k} = \sum_j u_{kj} x_{t,j} + b_k$, where $u_{kj}$ is the weight value associated with the $j$ th input unit and $k$ the output unit of the layer, $x_{t,j}$ and $b_j$ are the $j$ th elements of the input and bias vectors, it is easy to show that

$$\mathcal{R}(a_{t,k}) = \sum_j \mathcal{R}(u_{kj}) x_{t,j} + \sum_j u_{kj} \mathcal{R}(x_{t,j}) + \mathcal{R}(b_j)$$
$$= \sum_j v_{kj} x_{t,j} + \sum_j u_{kj} h'(a_{t,j}) \mathcal{R}(a_{t,j}) + v_j, \qquad (12)$$

where $v_{kj}$ and $v_j$ are the elements corresponding to $u_{kj}$ and $b_j$ of $\boldsymbol{v}$, $h(\cdot)$ is the hidden activation function transforming $a_{t,j}$ to $x_{t,j}$ by $x_{t,j} = h(a_{t,j})$, $a_{t,j}$ is an activation value produced by a previous layer. According to Eqn. (12), $\mathcal{R}(a_t^{\text{out}})$

can be calculated efficiently by modifying the forward propagation procedure, which results in the required directional derivative $\mathbf{J}\boldsymbol{v}$ since $\mathcal{R}(\boldsymbol{a}_t^{\text{out}}) = \nabla_\theta(\boldsymbol{a}_t^{\text{out}})\boldsymbol{v} = \mathbf{J}\boldsymbol{v}$. A detailed explanations of the modified forward procedure can be found in (Bishop, 2006).

In addition, it is not necessary to compute and store $\hat{\mathbf{H}}$ explicitly. For lattice-based MBR training, $\hat{\mathbf{H}}\,\mathcal{R}(\boldsymbol{a}_t^{\text{out}})$ can be directly calculated by

$$\hat{\mathbf{H}}\,\mathcal{R}(\boldsymbol{a}_t^{\text{out}}) = \kappa^2\gamma_t \odot \mathcal{R}(\boldsymbol{a}_t^{\text{out}}) - \kappa^2\gamma_t^{\text{MBR}}\left(\gamma_t^{\text{T}}\mathcal{R}(\boldsymbol{a}_t^{\text{out}})\right),$$

where $\odot$ refers to the Hadamard product.

## 4. Improving CG for Distributed Training

The CG-based distributed optimisation framework is presented in this section, which can be used for HF (Martens, 2010; Kingsbury et al., 2012). In practice, HF is found to have a high computational cost since a large number of CG iterations are required to perform in sequence to find an effective update (Sainath et al., 2013a). In this section, a modification to standard CG is presented that overcomes the numerical instability issue, which can reduce the number of CG iterations required by a factor of about twenty. Furthermore, a gradient normalisation method is proposed to improve the performance of CG for shared parameters, which facilitates the training of convolutional and recurrent models.

### 4.1. CG-based Distributed Optimisation

As shown in Fig. 1, the CG-based optimisation framework consists of two stages: the gradient accumulation stage and the CG stage. The gradient accumulation stage approximates the true gradient $\nabla_\theta \mathcal{L}(\theta)$ with an average of the gradients computed over every sample in a large data batch, which is referred to as a **gradient batch**. Here, a sample is an utterance for discriminative sequence training. In this stage, most of the calculations are used to compute the gradient w.r.t. each sample using the forward propagation and EBP procedures. These can be conducted in parallel using multiple workers. The negative gradients are then accumulated and averaged to form $\boldsymbol{v}_0$ for the CG stage.

In the CG stage, a sequence of CG iterations is used to find the parameter update $\Delta\theta$ with another batch of samples called the **CG batch**. As explained in Section 3.2, $\mathbf{G}\boldsymbol{v}_m$ is calculated at each CG iteration $m$ using the modified forward propagation and EBP. Such calculations can be conducted in parallel using a separate worker for each sample in the CG batch, whose output statistics can be accumulated for the rest of the steps of Alg. 1. The updated search direction $\boldsymbol{v}_{m+1}$ will monotonically improve upon $\boldsymbol{v}_m$ in terms of the loss value of the CG batch. After a certain number of iterations, $\Delta\theta_m$ with the best training loss performance on the CG mini-batch is returned as the direction found by the CG.

In this paper, executing the two stages once is referred to as an **update** since it is used to find one $\Delta\theta$ to update the current parameter $\theta$. In practice, we divide the whole training set randomly into $C$ partitions, with each of them being used as a gradient batch for an update, and hence each
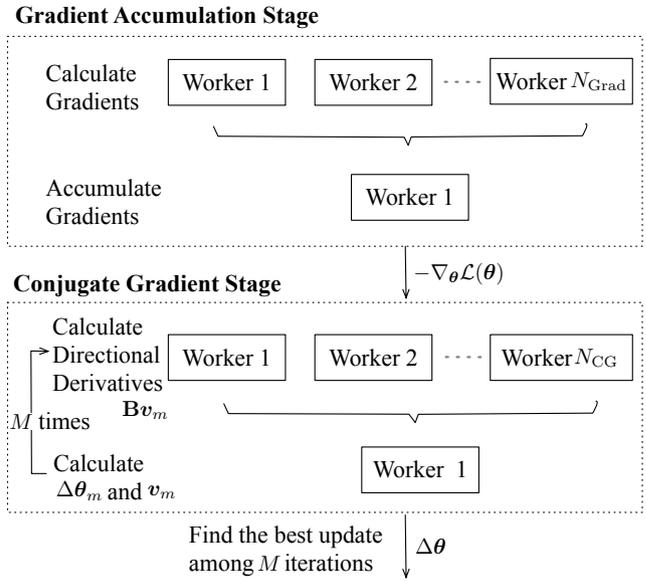
**Gradient Accumulation Stage**



**Figure 1**: A flow chart of our CG-based distributed optimisation framework, where $N_{\text{Grad}}$ and $N_{\text{CG}}$ are the utterance numbers in the gradient batch and CG batch respectively. They are also the maximum number of workers allowed for gradient calculation and directional derivative calculation. $\boldsymbol{v}_m$ and $\Delta\theta_m$ are the conjugate direction and the proposed parameter update at the $m$th CG iteration.

training epoch comprises $C$ updates performed in sequence. The CG batch is often much smaller than the gradient batch since it often needs to be processed for many iterations in the CG stage. In our experiments, we found it is better to sample the CG batch from the entire training set rather than just from the corresponding gradient batch.

### 4.2. Improving the Stability of CG

Although the calculation of an individual CG iteration can be distributed over many workers, CG iterations are still required to be performed sequentially. As reported in previous studies (Martens, 2010; Wiesler et al., 2013; Martens and Grosse, 2015), the HF method often requires about 200 CG iterations to find an effective update for training a DNN, even in the case of lattice-based discriminative sequence training (Kingsbury et al., 2012; Sainath et al., 2013a). This means that, in practice, CG restricts the training speed. Next, we explain the cause of this issue, and a solution is proposed to improve the stability of CG, which will be shown in Sec. 8 to yield effective updates only from few iterations of CG.

As discussed in Sec. 3.2, for a matching loss function, such as the CE loss together with the softmax function, the GN matrix $\mathbf{G}$ is in theory guaranteed to be positive semi-definite. However, even in our CE training experiments, it was observed that $\mathbf{G}$ could at times be negative (Haider, 2019). This issue was found to be a result of insufficient arithmetic accuracy when calculating the directional derivatives $\mathbf{J}\boldsymbol{v}_m$ for a CG iteration $m$. More specifically, let $\|\cdot\|_2$ be the L2 norm of a vector, when $\|\theta\|_2 \gg \|\boldsymbol{v}_m\|_2$, Eqn. (12)

becomes

$$\mathcal{R}(a_{t,k}) \approx \sum_j u_{kj} h'(a_{t,j}) \mathcal{R}(a_{t,j}),$$

due to the limited precision of the floating-point arithmetic, which can lead to an incorrect value of $\mathbf{G}$ that may no longer be a positive semi-definite matrix.

For large scale distributed training (Sainath et al., 2013a), the issue of negative $\mathbf{G}$ is resolved by using Tikhonov damping (Tikhonov et al., 1998), which uses $\mathbf{G} + \eta \mathbf{I}$ instead of $\mathbf{G}$ in CG. This corresponds to taking comparatively more conservative steps along the individual conjugate directions, and considerably slows down training as more CG updates are required to get a good overall solution. In the scenario when $\eta$ is large, Tikhonov damping is effectively analogous to an SGD step. Instead of Tikhonov damping, we propose to modify each CG iteration by using

$$\boldsymbol{v}'_m = (\|\boldsymbol{\theta}\|_2 / \|\boldsymbol{v}_m\|_2) \, \boldsymbol{v}_m$$

to compute $\mathbf{J}\boldsymbol{v}'_m$. Afterwards, $\mathbf{J}\boldsymbol{v}_m$ is obtained by

$$\mathbf{J}\boldsymbol{v}_m = (\|\boldsymbol{v}_m\|_2 / \|\boldsymbol{\theta}\|_2) \, \mathbf{J}\boldsymbol{v}'_m.$$

In our experiments, it was found that this improved CG algorithm can often produce an effective $\Delta\boldsymbol{\theta}$ with about 8 iterations.

### 4.3. Improving CG for Shared Parameters

In this section, it is demonstrated how CG can be adapted to perform efficiently for models with shared parameters, such as the TDNN (Waibel et al., 1989; Peddinti et al., 2015; Kreyssig et al., 2018) and LSTM (Hochreiter and Schmidhuber, 1997; Graves et al., 2013; Sak et al., 2014), both widely used for acoustic modelling. In contrast to a DNN, a TDNN uses a sequence of fully FC layers to perform 1-dimensional (-dim) convolutions across time, whose input vectors the concatenation of $\boldsymbol{x}_{t_1}$ and $\boldsymbol{x}_{t_2}$, the output from their direct preceding layers of two different time steps $t_1$ and $t_2$. That is,

$$\boldsymbol{y}_t = h(\mathbf{U}\,\text{Concat}(\boldsymbol{x}_{t_1}, \boldsymbol{x}_{t_2}) + \boldsymbol{b}),$$

where $\text{Concat}(\cdot)$ is the concatenation operation; $h(\cdot)$, $\mathbf{U}$ and $\boldsymbol{b}$ are the activation function, weight matrix and bias vector of the FC layer. Hence, the directional derivatives of a TDNN are also calculated using Eqn. (12). Alternatively, a TDNN can be viewed as a feedforward model with a binary tree structure by duplicating each FC layer for the relevant time steps, and the parameters $\mathbf{U}$ and $\boldsymbol{b}$ are shared across time.

Regarding an (Elman network) RNN (Rumelhart et al., 1986; Elman, 1990), $\boldsymbol{h}_t$, the output vector at time $t$, is generated by transforming a concatenation of $\boldsymbol{h}_{t-1}$ and the current input $\boldsymbol{x}_t$ with an FC layer whose parameters are $\mathbf{U}$ and $\boldsymbol{b}$, i.e. $\boldsymbol{y}_t = f(\mathbf{U}\,\text{Concat}(\boldsymbol{x}_t, \boldsymbol{h}_{t-1}) + \boldsymbol{b})$. An LSTM is an improved RNN with enhanced long-term memory capability based on the gating mechanism, which uses $\boldsymbol{i}_t$, $\boldsymbol{f}_t$, and $\boldsymbol{o}_t$, the sigmoidal output vectors from three extra FC layers, to simulate the logic gates of a memory circuit to maintain a memory cell $\boldsymbol{c}_t$. Given the parameters of the three extra FC

layers, $\mathbf{U}_i$ and $\boldsymbol{b}_i$, $\mathbf{U}_f$ and $\boldsymbol{b}_f$, and $\mathbf{U}_o$ and $\boldsymbol{b}_o$, an LSTM layer can be specifically presented as

$$
\begin{aligned}
\boldsymbol{i}_t &= \sigma(\mathbf{U}_i\,\text{Concat}(\boldsymbol{x}_t, \boldsymbol{h}_{t-1}) + \boldsymbol{b}_i), \\
\boldsymbol{f}_t &= \sigma(\mathbf{U}_f\,\text{Concat}(\boldsymbol{x}_t, \boldsymbol{h}_{t-1}) + \boldsymbol{b}_f), \\
\boldsymbol{c}_t &= \boldsymbol{f}_t \odot \boldsymbol{c}_{t-1} + \boldsymbol{i}_t \odot \tanh(\mathbf{U}\,\text{Concat}(\boldsymbol{x}_t, \boldsymbol{h}_{t-1}) + \boldsymbol{b}), \\
\boldsymbol{o}_t &= \sigma(\mathbf{U}_o\,\text{Concat}(\boldsymbol{x}_t, \boldsymbol{h}_{t-1}) + \boldsymbol{b}_o), \\
\boldsymbol{h}_t &= \boldsymbol{o}_t \odot \tanh(\boldsymbol{c}_t),
\end{aligned}
$$

where $\sigma(\cdot)$ and $\tanh(\cdot)$ are the sigmoid and hyperbolic tangent activation functions. An LSTM layer can be implemented with four FC layers and the Hadamard product. Therefore, the directional derivative for an LSTM can be calculated using Eqn. (12) and the following rule of the $\mathcal{R}(\cdot)$ operator for gating:

$$\mathcal{R}(\boldsymbol{g}_t \odot \boldsymbol{z}_t) = \mathcal{R}(\boldsymbol{g}_t) \odot \boldsymbol{z}_t + \boldsymbol{g}_t \odot \mathcal{R}(\boldsymbol{z}_t), \tag{13}$$

where $\boldsymbol{g}_t$ and $\boldsymbol{z}_t$ are two example vectors. Further by unfolding through time for $u$ steps (Robinson and Fallside, 1987; Werbos, 1988), a folded LSTM layer becomes $u$ unfolded layers with feedforward connections sharing all of their parameters. The input of the $v$ th feedforward layer is a concatenation of the output from the $(v-1)$ th layer and $\boldsymbol{x}_{t-u+v}$.

Next we present a modification to the CG algorithm to make the algorithm more effective for models with shared parameters. From Alg. 1, both the step size $\alpha_m$ and the conjugate search direction $\boldsymbol{v}_m$ are determined by the dot product of the residual $\boldsymbol{r}_m^{\mathrm{T}} \boldsymbol{r}_m$ and the directional derivative $\mathbf{G}\boldsymbol{v}_m$. For models such as the TDNN and RNN, shared parameters receive more updates and hence will contribute more to the norm of the vectors $\boldsymbol{r}_0$ and $\mathbf{G}\boldsymbol{v}_m$ than the parameters which are not shared. Careful preliminary experiments using TDNNs and LSTMs found that in situations where the shared parameters dominate the norm of these vectors, the CG algorithm was slow to find an update direction that could reduce the loss value. Our solution is apply a diagonal scaling to $\boldsymbol{r}_m$ and $\mathbf{G}\boldsymbol{v}_m$ by a matrix whose diagonal entries correspond to the square root of the number of times that a parameter is shared. This ensures that the L2 norm of the vectors are not dominated by the contributions of the shared parameters and enables a more effective update for the other parameters. Such an approach corresponds to preconditioning the CG algorithm (Shewchuk, 1994) by applying the diagonal scaling only to $\boldsymbol{r}_0$ among all the residuals $\boldsymbol{r}_m$.

By using the algorithm given in (Zhang and Woodland, 2015), the diagonal scaling can be efficiently achieved by adding an operation to normalise the resulting gradients or directional derivatives at the end of the EBP procedure by the number times a parameter is shared. Experiments showed that this solution enabled CG to find progressively better update directions in each CG iteration for TDNNs and LSTMs (Haider, 2019). More details of the theoretical analysis and experimental evidence can be found in (Haider, 2019).

## 5. Natural Gradient Optimisation

In this section first NG is reviewed. Then the proposed CG-based optimisation framework for NG for discriminative

sequence training is presented.

## 5.1. Natural Gradient Descent

This section presents an overview of NG descent (Amari, 1997; Pascanu et al., 2013). Let $\mathcal{M}$ denote the space of all temperature-modulated probability distributions $P_\theta(\mathbf{W}|\mathbf{O})^4$. In the context of optimisation, NG tries to find an update $\Delta\theta$ that minimises the loss function locally while keeping a similar probability distribution to the one resulting from the current $\theta$. That is,

$$\Delta\hat{\theta} = \arg\min_{\Delta\theta} \mathcal{L}(\theta + \Delta\theta), \quad (14)$$

$$\text{s.t. } \mathbb{E}_{p(\mathbf{O})} \left[ \text{KL}\left( P_\theta(\mathbf{W}|\mathbf{O}) \| P_{\theta+\Delta\theta}(\mathbf{W}|\mathbf{O}) \right) \right] \leqslant \epsilon$$

where $\text{KL}(\cdot\|\cdot)$ refers to the KL-divergence and $\epsilon$ is a constant controlling the speed of exploration along the manifold. In Eqn. (14), we assume $p_\theta(\mathbf{W}, \mathbf{O}) = P_\theta(\mathbf{W}|\mathbf{O})p(\mathbf{O})$.

Eqn. (14) can be formulated as an equivalent constrained optimisation problem in the parameter space by using *Lagrange multipliers*, which restrict the exploration of the parameter space to be within a local neighbourhood of the current estimate $\theta$. Approximating $\mathcal{L}(\theta + \Delta\theta)$ by its first-order Taylor approximation $\mathcal{L}(\theta) + \Delta\theta^T\nabla_\theta\mathcal{L}(\theta)$, and the KL-divergence constraint by its second-order Taylor approximation (Amari, 2016) yields

$$\Delta\hat{\theta} = \arg\min_{\Delta\theta} \left\{ \mathcal{L}(\theta) + \Delta\theta^T\nabla_\theta\mathcal{L}(\theta) + \frac{\lambda}{2}\Delta\theta^T\mathbf{F}\Delta\theta \right\}, \quad (15)$$

where $\lambda$ is the Lagrange multiplier that controls the compromise between minimising the loss and satisfying the KL-divergence constraint. $\mathbf{F}$ is referred to as the FI matrix.

$$\mathbf{F} = \mathbb{E}_{p_\theta(\mathbf{W}|\mathbf{O})}\left[ \nabla_\theta \ln P_\theta(\mathbf{W}|\mathbf{O})\nabla_\theta \ln P_\theta(\mathbf{W}|\mathbf{O})^T \right], \quad (16)$$

Similar to second-order methods, NG attempts to minimise a quadratic function at each iteration. Differentiating Eqn. (15) and equating it to zero yields the solution

$$\Delta\theta = -\frac{1}{\lambda}\mathbf{F}^{-1}\nabla_\theta\mathcal{L}(\theta). \quad (17)$$

This suggests that the update direction obtained by NG transforms the steepest decent direction by taking into account the curvature information of the log-likelihood given by $\mathbf{F}^{-1}$. It can be shown that such a direction is indeed the optimal descent in the loss surface generated on the manifold $\mathcal{M}$.

Computing the exact FI matrix defined in Eqn. (16) requires taking the expectation over the distribution $p_\theta(\mathbf{W}, \mathbf{O})$ which is infeasible for LVCSR. A standard approach is to approximate this expectation by using the average of samples from $p_\theta(\mathbf{W}, \mathbf{O})$. In this paper, we take samples from the data distribution $p(\mathbf{W}, \mathbf{O})$, instead of $p_\theta(\mathbf{W}, \mathbf{O})$. This form of the FI matrix yields the empirical Fisher matrix where the contribution of each utterance corresponds to

$$\mathbf{F} \approx \nabla_\theta \ln P_\theta(\mathbf{W}^{\text{ref}}|\mathbf{O})\nabla_\theta \ln P_\theta(\mathbf{W}^{\text{ref}}|\mathbf{O})^T. \quad (18)$$

---

[4] $P_\theta(\mathbf{W}|\mathbf{O})$ is determined by a chosen value of $\kappa$ (see Eqn. (2)).

As reviewed in Section 2.2, to reduce the difficulty in calculating $\mathbf{F}^{-1}$, $\mathbf{F}$ is often assumed to have a (block) diagonal structure (Povey et al., 2015; Kingma and Ba, 2015; Martens and Grosse, 2015) with the diagonal blocks corresponding to Kronecker products of two smaller matrices. Meanwhile, CG was proposed to compute the NG direction without calculating $\mathbf{F}^{-1}$ explicitly for both reinforcement learning (Schulman et al., 2015) and discriminative sequence training (Haider and Woodland, 2017). This will be presented in detail later in Sec. 5.2.

There exist theoretical advantages in using NG. For ML training, assuming the distribution of the gradient of the expected loss to be an isotopic Gaussian, the NG direction is relevant to a direction in the parameter space that maximises the probability of reducing the generalisation error (Roux et al., 2008). Recently, it has been shown that for convex problems, solving Eqn. (15) at each iteration results in an exponentially faster convergence speed compared to gradient descent (Bernacchia et al., 2018).

## 5.2. Natural Gradient with CG for Discriminative Sequence Training

From Eqn. (18), the empirical FI matrix $\mathbf{F}$ is guaranteed to be positive semi-definite and therefore the optimisation framework proposed in Sec. 4 can be used for NG. CG is used to solve $\lambda\mathbf{F}\Delta\theta = -\nabla_\theta\mathcal{L}(\theta)$ (Eqn. (17)). From Eqn. (2), $\ln P_\theta(\mathbf{W}^{\text{ref}}|\mathbf{O}) = -\mathcal{L}_{\text{MMI}}(\theta)$ when $\kappa = 1$, and hence Eqn. (18) can be re-written as

$$\mathbf{F} \approx \nabla_\theta\mathcal{L}_{\text{MMI}}(\theta)\nabla_\theta\mathcal{L}_{\text{MMI}}(\theta)^T$$

$$= \sum_{t=1}^{T} \mathbf{J}^T\nabla_{a_t^{\text{out}}}\mathcal{L}_{\text{MMI}}(\theta)\nabla_{a_t^{\text{out}}}\mathcal{L}_{\text{MMI}}(\theta)^T\mathbf{J}. \quad (19)$$

This uses the fact that $\nabla_\theta\mathcal{L}_{\text{MMI}}(\theta) = [\mathbf{J}^T\nabla_{a_t^{\text{out}}}\mathcal{L}_{\text{MMI}}(\theta)]_{t=1}^{T}$, where $\mathbf{J}$ is the Jacobian matrix of $a_t^{\text{out}}$ w.r.t. $\theta$, and $T$ is the number of frames in the utterance. Thereafter, we denote $\hat{\mathbf{F}}$ as $\nabla_{a_t^{\text{out}}}\mathcal{L}_{\text{MMI}}(\theta)\nabla_{a_t^{\text{out}}}\mathcal{L}_{\text{MMI}}(\theta)^T$ for simplicity. Since Eqn. (19) has the same form as the GN matrix given in Eqn. (10), the procedure given in Sec. 3.4 can be used to calculate $\mathbf{F}v$ (within a CG iteration), which first calculates $\mathbf{J}v$ using the modified forward propagation, then multiplies the resulting vector by $\hat{\mathbf{F}}$, and at last calculates $\mathbf{J}^T(\hat{\mathbf{F}}\mathbf{J}v)$ using the EBP procedure.

Next, it is shown how $\hat{\mathbf{F}}$ can be calculated. Similar to ML and MBR training discussed in Sec. 3.4, it is easy to show $\partial\mathcal{L}_{\text{MMI}}(\theta)/\partial a_{t,i} = -\kappa\gamma_{t,k}^{\text{MMI}}$ (Zhang, 2017), where $a_{t,k}$ is the logit value at time $t$ of output unit $k$, and $\gamma_{t,k}^{\text{MMI}}$ is the MMI occupancy with $\gamma_{t,k}^{\text{MMI}} = \gamma_{t,k}^{\text{num}} - \gamma_{t,k}^{\text{den}}$. $\gamma_{t,k}^{\text{num}}$ and $\gamma_{t,k}^{\text{den}}$ are the occupancy derived separately from the numerator and denominator parts of Eqn. (2).

In practice, $\hat{\mathbf{F}} = \kappa^2\gamma_t^{\text{MMI}}(\gamma_t^{\text{MMI}})^T$ does not need to be calculated and stored explicitly. Recall the $\mathcal{R}(\cdot)$ operator discussed in Sec. 3.2, and that $\mathcal{R}(a_t^{\text{out}}) = \mathbf{J}v$, which has the same dimension as $\gamma_t^{\text{MMI}}$. To calculate $\hat{\mathbf{F}}\mathbf{J}v$ directly, $(\gamma_t^{\text{MMI}})^T\mathcal{R}(a_t^{\text{out}})$ can be obtained first, which is followed by scaling $\gamma_t^{\text{MMI}}$. Specifically,

$$\hat{\mathbf{F}}\mathbf{J}v = \kappa^2\gamma_t^{\text{MMI}}\left( (\gamma_t^{\text{MMI}})^T\mathcal{R}(a_t^{\text{out}}) \right).$$

The term $(\gamma_t^{\text{MMI}})^{\text{T}}\mathcal{R}(\boldsymbol{a}_t^{\text{out}})$ is a scalar quantity and can be interpreted as a learning rate for $\gamma_t^{\text{MMI}}$.

By comparing Eqns. (4) and (15), the difference between the NG and HF approach when applied to minimise any arbitrary smooth loss function lies in the matrix used in the second-order term, $\mathbf{G}$ and $\mathbf{F}$. For HF, the GN matrix $\mathbf{G}$ requires calculating the Hessian of the MBR loss w.r.t. the logit values. For NG, the empirical FI matrix is calculated as the outer product of the gradients of the MMI loss w.r.t. the logit values, regardless of the loss used for training. In the case of MBR training, NG provides an efficient procedure to combine MBR loss with MMI loss, which is similar to the widely used "MMI prior" method for GMM-HMM MBR training that interpolates MBR with MMI in the loss function (Young et al., 2015; Zhang and Woodland, 2017).

## 6. Regulating NG Updates

In Sec. 5 it is shown how NG descent corresponds to scaling and rotating the gradient $\nabla_\theta \mathcal{L}(\theta)$ through multiplication with the inverse $\mathbf{F}$ matrix. In a scenario where the training criterion is only an approximation of the evaluation metric, it is shown in (Haider and Woodland, 2018; Haider, 2019) that both the approximate NG descent and SGD can at times follow a path in the parameter space where generalisation improvements w.r.t. the training criterion on the validation set fail to correlate with reductions in WER (the evaluation criterion of interest). Thus in such training paradigms, over-fitting can occur not only due to the lack of training data but also due to the underlying criterion mismatch. Depending on the task, it will be attractive to have a mechanism that regulates the amount of scaling and rotation of the loss gradient to achieve good generalisation. The following subsection presents a common framework to regulate the NG direction or the gradient descent direction by using an appropriate choice of $\mathbf{B}$. The procedure presented here relies on a re-derivation of Taylor's theorem using the concepts of manifolds, tangent vectors and directional derivatives from the perspective of differential geometry. The derivation is provided in detail in the technical report (Haider, 2018). An overview of the necessary underlying concepts can be found in (Amari, 2016; de Felice and Clarke, 1992).

### 6.1. A Common Framework to Regulate Natural Gradient and Gradient Descent

Assuming that the loss function $\mathcal{L}(\theta)$ is sufficiently smooth, using Taylor's theorem, second order methods proceed to minimise the loss by minimising the following function at each iteration $m$:

$$\Delta\theta' = \underset{\Delta\theta}{\arg\min}\left\{\mathcal{L}(\theta_m) + \Delta\theta^{\text{T}}\nabla_\theta\mathcal{L}(\theta_m) + \frac{1}{2}\Delta\theta^{\text{T}}\mathbf{B}\Delta\theta\right\}$$

where $\theta_m$ corresponds to the current parameter estimate. Using the fundamental theorem of calculus, in (Haider, 2018), it is shown that solving the above problem can be cast as an equivalent minimisation problem in the tangent space of the

current parameter estimate $\mathcal{T}(\theta_m)$:

$$\underset{\Delta\theta\in\mathcal{T}(\theta_m)}{\arg\min}\left\{\mathcal{L}(\theta_m) + \langle\Delta\theta^{\text{T}}, \nabla_\theta\mathcal{L}(\theta_m)\rangle + \frac{1}{2}\Delta\theta^{\text{T}}\mathbf{B}\Delta\theta\right\}.$$

To aid understanding, the notion of the tangent space associated with a point $\theta_m$ in the parameter space corresponds to the set of all possible directions that can be traversed from $\theta$ that yields different directional derivatives w.r.t. the loss function $\mathcal{L}$, a map from the parameter space to the real line. Thus, the $\Delta\theta$ that are often probed in the optimisation are in fact members of $\mathcal{T}(\theta_m)$. In the above equation, $\langle\Delta\theta^{\text{T}}, \nabla_\theta\mathcal{L}(\theta_m)\rangle$ corresponds to an inner product between vectors in $\mathcal{T}(\theta_m)$ where the inner product can be generalised to any Riemannian metric (de Felice and Clarke, 1992). Replacing the standard inner product with the positive definite $\mathbf{F}^{-1}$ leads[5] to the following optimisation problem in $\mathcal{T}(\theta_m)$:

$$\underset{\Delta\theta\in\mathcal{T}(\theta_m)}{\arg\min}\left\{\mathcal{L}(\theta_m) + \Delta\theta^{\text{T}}\mathbf{F}^{-1}\nabla_\theta\mathcal{L}(\theta_m) + \frac{1}{2}\Delta\theta^{\text{T}}\mathbf{B}\Delta\theta\right\}$$
(20)

where each entry of the matrix $\mathbf{F}^{-1}$ is a smooth function of the current estimate $\theta_m$. Differentiating Eqn. (20) and equating to zero leads to following solution:

$$\mathbf{B}\Delta\theta = -\mathbf{F}^{-1}\nabla_\theta\mathcal{L}(\theta).$$
(21)

The right hand side of this above equation corresponds to the NG direction. Hence Eqn. (20) presents a procedure to regulate the scaling and rotation applied by $\mathbf{F}^{-1}$ with a matrix $\mathbf{B}$ chosen in an appropriately understood sense.

### 6.2. NGHF: Combining NG and HF based on CG

The choice of appropriate $\mathbf{B}$ varies from task to task. In (Haider and Woodland, 2018), the prevalence of over-fitting due to criterion mismatch was observed to be highly correlated with the increased sharpness of the NN frame posteriors. In Sec. 3.3, it is discussed how scaling with the inverse of the GN matrix regulates sharp changes in the entropy of DNN frame posteriors. Using this insight, in this work $\mathbf{G}$ is employed to regulate the NG descent direction. Computing the individual inverse matrix scalings in Eqn. (20) directly is expensive in terms of both computation and storage. Using the HF approach, each individual matrix scaling is approximated by solving equivalent linear systems using CG. In Sec. 3 it is shown how the update direction proposed at each CG iteration corresponds to:

$$\Delta\theta_{k+1} \leftarrow \Delta\theta_k + a_k\boldsymbol{v}_k$$

where $\boldsymbol{v}_k$ represents the current conjugate direction. At the first iteration of CG, this is the direction that the algorithm has been initialised with. In contrast to NG and HF, the initial direction now corresponds to approximation of the NG direction instead of the gradient. Thus when Eqn. (20) is solved with CG, the resultant update corresponds to:

$$\Delta\theta = w_1\Delta\theta_{\text{NG}} + w_2\Delta\theta_{\text{HF}}$$
(22)

---

[5]The derivation holds for any positive scalar multiple of the Fisher.

which is a weighted combination of the NG direction and conjugate directions computed using local curvature information. Hence, in this sense we denote this approach NGHF.

# 7. Experimental Setup

The proposed optimisation framework was evaluated on data from the multi-genre broadcast (MGB) challenge (Bell et al., 2015) which uses data from a wide range of BBC television programmes, and the effectiveness of the techniques for discriminative sequence training with the MPE loss was found. All systems were trained using a 200 hour training set. The official development set dev.full was split into two subsets. One split corresponds to the official MGB subset, **dev.sub**, with 5.5 hours data, which is used as the validation set to choose the hyper-parameters and select the best parameter update $\Delta\theta$ in Alg. 1. To evaluate the generalisation ability to unseen data, an evaluation test set **dev.sub2** was also created, which consists of 23 hours data from the remaining 35 episodes in the dev.full set. Further detail related to the data preparation can be found in (Woodland et al., 2015). The input to all models were 40-dim log-Mel filter bank features extended with their delta coefficients, which were normalised at the utterance-level for mean and at the show-level for variance (Woodland et al., 2015). All experiments were conducted using HTK version 3.5 and the Py-HTK pipelines(Young et al., 2015; Zhang et al., 2019a).

The RNN models used in the experiments consist of two 1000-dim recurrent layers followed by a 1000-dim feedforward layer. Each recurrent layer is unfolded for 20 steps (from +5 to −14). Apart from replacing the standard RNN layer with the LSTM layers of the same size, the LSTMs have the same structure as the RNNs. The TDNNs used the same structure as in (Peddinti et al., 2015) with five 1000-dim hidden layers, whose context shifts used to splice the features are $\{-2, -1, 0, 1, 2\}$, $\{-1, 2\}$, $\{-3, 3\}$, $\{-7, 2\}$ and $\{0\}$ from the input to the output layers respectively. For all models, the output layer consists of about six thousand output units, with each output corresponding to a context-dependent triphone state obtained by conventional decision tree tying approach.

The large-batch-based methods, HF, NG, and NGHF, are compared with the mini-batch-based methods, SGD and Adam. For this the LSTM-HMMs, RNN-HMMs, and TDNN-HMMs are trained at the sequence-level with these different optimisers based on the MPE loss. Prior to sequence training, they are trained with SGD on a frame-level CE loss. To monitor the occurrence of over-fitting and avoid the mismatch between the WER and MPE loss, decoding was performed on both of the validation and evaluation sets, using a 158k word vocabulary trigram language model.

Since standard SGD and Adam suffer from high data transmission costs in synchronous distributed processing, such optimisers used single GPU training and standard configurations (Su et al., 2013; Veselỳ et al., 2013; Zhang and Woodland, 2015). The hyper-parameters associated with SGD and Adam were chosen using grid search such that the

improvements obtained from MPE training are closely correlated with WER reductions on the validation set.

Following (Kingsbury et al., 2012), all experiments on HF, NG, and NGHF that used the proposed CG-based optimisation framework were performed in a synchronous distributed setting where the gradients were computed across multiple workers in parallel and then accumulated. To achieve a good balance between the reduced training time cost through parallelisation and keeping the data transmission cost low, the gradient batch used in this work contains 25 hours of data. When using four workers instead of one worker to collect the gradients over a 25 hour gradient batch, the time cost is reduced from 150 minutes to 37 minutes, which is an almost a linear reduction in the time cost by a factor of the number of workers. The method itself allows full parallelisation of the gradient calculation stage such that one worker is used per utterance. Parallelisation is not only constrained by the number of workers available, but also the additional overhead time-cost due to further parallelisation and data transmission[6]. To reduce the time cost associated with the individual CG iterations, a 0.5 hour CG batch was used throughout the experiments, whose data were uniformly sampled from the entire training set. It should be noted that one worker was used for CG, although in theory the maximum number of workers allowed for CG could be the same as the number of utterances in the CG batch. It was found that running CG for 5-8 iterations was sufficient to find an parameter update that could yield a reasonable reduction in the loss value. In terms of computational cost, the CG worker was found to take about 30 minutes to execute 8 CG iterations for LSTM-HMMs. All timing information was obtained on a machine with four Tesla P100 GPUs and a 14 core Intel Xeon CPU E5-2680 v4 at 2.40GHz.

From Table 1, it is clear that validating the performance at each CG iteration takes the largest proportion of the overall time cost. Recalling that the goal of CG is to minimise a quadratic approximation of the loss function, through the validation of the update obtained by each CG iteration on the CG batch, the validation stage checks whether CG follows a path in the parameter space where the quadratic approximation still holds. Although a validation stage was performed after every CG iteration in all the experiments in this paper, we empirically found that the check can be performed less frequently to reduce the time cost.

# 8. Experimental Results

This section first compares the various optimisers using LSTM models (Sec. 8.1). The evolution of the MPE loss function during discriminative sequence training is discussed, as well as the final WER on the evaluation set. The performance when using both sigmoid and ReLU activation functions with RNN and TDNN models is then investigated for the different optimisation approaches (Sec. 8.2).

---

[6]The optimal degree of parallelisation will vary depending on the computing infrastructure used.

**Table 1**
The proportion of time cost for running the CG algorithm for NGHF for 8 iterations using a single NVidia P100 GPU and a single core of Intel Xeon CPU E5-2680 v4 at 2.40GHz. The time cost for loading the lattices into the memory is excluded from the calculation since it can vary considerably depending on the lattice implementation and the speed of the storage system.

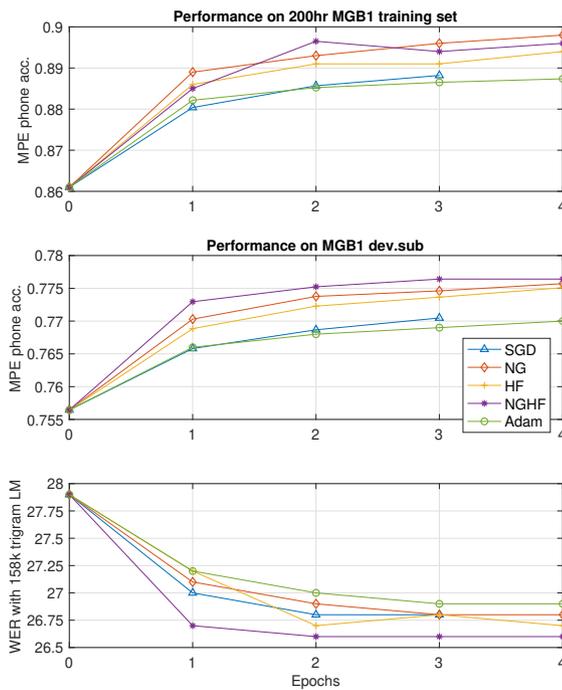| Procedure | %Time cost |
|---|---|
| Modified forward propagation | 15.1 |
| EBP | 7.8 |
| Collecting statistics over lattices | 4.1 |
| Evaluating the performance of each $\Delta\theta$ | 73.0 |



**Figure 2:** The evolution of the performance of the LSTM-HMMs with different optimisers. The first two plots show the phone accuracy on a subset of the training set and the validation set dev.sub. The third plot shows the absolute reductions in %WERs on the validation set.

## 8.1. LSTM-HMMs Results

Figure 2 shows how MPE phone accuracy (negative of the MPE loss) evolves during training for the LSTM-HMMs with different optimisers, and the corresponding changes in WER. Table 2 summarises the results on the validation set for the best epoch. From Fig. 2, NGHF is the most effective method in terms of both improving the MPE accuracy and reducing the WER. By combining the KL-divergence with the local curvature information obtained from the Hessian matrix, the NGHF method achieves a greater WER reduction with just 8 parameter updates (*i.e.* one epoch) than the

**Table 2**
LSTM-HMM performance on the validation set with different optimisers. "Best epoch" gives the epoch model with the best performance. "#Update" shows the number of updates used in the MPE training, and "k" stands for one thousand. "MPE Acc." gives the MPE accuracy (negative of MPE loss value). Although NG, HF, and NGHF used much fewer updates, the numbers of utterances processed in each epoch are similar to those of SGD and Adam.

| Optimiser | Best epoch | #Update | MPE Acc | %WER |
|---|---|---|---|---|
| CE | — | — | 0.765 | 27.9 |
| SGD | 3 | 420k | 0.771 | 26.8 |
| Adam | 4 | 560k | 0.770 | 26.9 |
| NG | 3 | 24 | 0.775 | 26.8 |
| HF | 2 | 16 | 0.772 | 26.7 |
| NGHF | 2 | 16 | 0.775 | **26.6** |

**Table 3**
%WERs for MPE trained LSTM-HMMs on the evaluation set with different optimisers. "CE" refers to the SGD trained CE system, while the others refer to the MPE systems trained with the corresponding optimisers.

| CE | SGD | Adam | NG | HF | NGHF |
|---|---|---|---|---|---|
| 29.3 | 28.6 | 28.6 | 28.6 | 28.4 | **28.3** |

converged model obtained trained using SGD or Adam using hundreds of thousands of updates. Table 2, shows that MPE training using NGHF results in a relative 4.7% WER reduction (WERR) over the CE trained model. It can also be seen that all of the three CG-based second-order optimisers require far fewer updates and fewer training epochs to converge to a good solution than SGD and Adam.

To ensure the performance difference of LSTM-HMMs trained with different optimisers can be generalised to unseen data, all models were tested on the evaluation set, which was not use for setting hyper-parameters. Table 3 summarises the results of the LSTM-HMMs on the evaluation set. From the results in Table 3, MPE with NGHF achieves the lowest WERs among all five optimisers, and obtains a 3.4% relative WERR compared to the CE trained model. Compared to both SGD and Adam, NGHF achieves a 1% larger WERR while requiring far fewer parameter updates and also fewer training epochs. This improvement was found to be statistically significant at the 0.1% level for dev.sub2[7].

## 8.2. TDNN-HMMs and RNN-HMMs with sigmoid and ReLU Activation Functions

This section presents a comparison between RNNs and TDNNs with both ReLU and sigmoid activation functions when using the various optimisers for MPE training. It is well-known that ReLU and sigmoid models require very different configurations in SGD-based training: ReLU models

---

[7]Statistical significance tests in this paper use a sign test on WER differences at the episode level for the 35 programme episodes included in dev.sub2.

**Table 4**

Number of MPE updates required by different optimisers ("k" stands for 1000) for ReLU and sigmoid ($\sigma$) RNN and TDNN models. Although NG, HF, and NGHF used far fewer updates, the numbers of utterances processed in each epoch are similar to those of SGD and Adam.

| Model | SGD | Adam | NG | HF | NGHF |
|---|---|---|---|---|---|
| ReLU RNN | 100k | 100k | 0 | 24 | 16 |
| $\sigma$ RNN | 440k | 440k | 40 | 40 | 40 |
| ReLU TDNN | 330k | 100k | 0 | 32 | 32 |
| $\sigma$ TDNN | 440k | 440k | 48 | 48 | 40 |

**Table 5**

%WER for ReLU and sigmoid ($\sigma$) RNN and TDNN models on the evaluation set with different optimisers. "CE" refers to the SGD trained CE system, while other columns refer to the MPE systems trained with the corresponding optimisers.

| Model | CE | SGD | Adam | NG | HF | NGHF |
|---|---|---|---|---|---|---|
| ReLU RNN | 30.3 | 30.3 | 30.2 | 30.3 | **29.6** | 29.7 |
| $\sigma$ RNN | 32.2 | 31.6 | 31.6 | 30.6 | 30.7 | **30.5** |
| ReLU TDNN | 30.6 | 29.8 | 29.8 | 30.6 | 29.6 | **29.3** |
| $\sigma$ TDNN | 29.9 | 28.5 | 28.3 | 28.2 | 28.5 | **28.1** |

often need a learning rate and the range for random initialisation a factor of 4 to 8 times smaller than those used for sigmoid models, and that discriminative sequence training with ReLU activation functions often results in over-fitting to MBR loss functions[8]. We claim that such a difference still exists with the second-order optimisers since the second-order derivatives of ReLU and sigmoid are very different[9], which causes a large difference in the calculation of directional derivatives (Bishop, 2006).

Table 4 gives the number of updates required by each type of optimiser to find a good solution and Table 5 compares the efficacy of the trained models with different optimisers on the evaluation set. For the sigmoid models, using NGHF obtains the largest reductions in WER from MPE training, yielding a 5% WERR for RNN and 6% WERR for TDNN, over the relevant CE trained models. Compared to SGD and Adam, NGHF achieves a WERR of 3.5% for sigmoid RNN and 1% for sigmoid TDNN while using far fewer updates. These WER reductions were found to be statistically significant (p < 0.001) on dev.sub2.

Regarding the ReLU models, it was found to be very difficult to get reasonable WER reductions using SGD, Adam and NG. With NG in particular, sequence training was found to suffer from over-fitting due to the mismatch between the MPE loss and the WERs from the very beginning of the MPE training. On the validation set, NG was observed to achieve excellent generalisation performance in terms of MPE loss but such improvements failed to correlate well with the WER reductions. For the TDNN ReLU model, MPE training with NGHF gives the the biggest WER reduction, which achieves a WERR of 4.2% compared to the CE trained model. In contrast to SGD and Adam, NGHF achieves WERRs of 2%. For the ReLU RNN model, although HF produced the lowest WER, its WERR over the NGHF method was not found to be statistically significant (at the 5% level).

## 9. Discussion

The optimisation framework presented in this work provides a general method to flexibly train NN models with

any arbitrary smooth loss using a large batch NG descent or second-order method in a data-parallel manner. Although the work has applied the framework in a centralised distributed training environment, the framework presented can also be applied in a decentralised distributed training environment as well (Zhang et al., 2019b,c).

Compared to the traditional HF approach, the novel modifications proposed in this paper overcome the computational drawbacks of requiring hundreds of CG iterations that have been cited as a key issue with HF (Martens and Grosse, 2015). The experimental results presented with the LSTM, TDNN and RNN models show how effective updates from only a small number of CG iterations can be obtained, and this approach to model training leads to significant WER reductions. The efficacy of the novel preconditioning approach mentioned in Sec. 4.3 can be clearly seen when training the LSTM model. Table 3 shows that NG, HF and NGHF all yield greater WER reductions in comparison to the stochastic-gradient-based approaches while using far fewer updates. Furthermore, this paper presents a novel procedure to regulate NG learning whose efficacy can be seen in Table 4 and 5. For the ReLU-based TDNN and RNN models, the regularised NG approach (*i.e.* NGHF) is more adept in following a path in the parameter space where minimising w.r.t. the training criterion correlates with WER reductions. This can be very important when applying NG to other structures and non-ASR tasks, since ReLU is widely used by the models of computer vision (Simonyan and Zisserman, 2015) and natural language processing (Vaswani et al., 2017).

In addition to the differences in WER and the number of updates reported in Sec. 8, each of the optimisers require different amounts of computation and GPU memory. For SGD, the input and output values of each layer are calculated in forward-propagation. In the backward-propagation pass, the derivatives w.r.t. the input/output values, as well as the gradients of the parameters are calculated. The parameters $\theta$, gradients $\nabla_{\theta}\mathcal{L}(\theta)$, and update values $\Delta\theta$ each require the same amount of space when they are stored in GPU memory. The input/output values and their derivatives are also stored in GPU memory, and the required storage depends on the number of frames in the mini-batch. For Adam, extra calculation and GPU memory are required to compute and store $(\nabla_{\theta}\mathcal{L}(\theta))^2$ and the relevant second raw moment estimate for every frame in the training set (Kingma and Ba, 2015).

NG, HF, and NGHF all have the same computation and

---

[8]By over-fitting we here mean a low correlation between reduction in MBR loss and reduction in WER.

[9]For ReLU $h''(a_{t,j}) = 0$ and for sigmoid $\sigma''(a_{t,j}) = \sigma(a_{t,j})(1 - \sigma(a_{t,j}))(1 - 2\sigma(a_{t,j}))$

storage complexity as SGD for the gradient accumulation stage. In addition to the standard forward-propagation and backward-propagation procedures, each CG iteration uses an extra modified forward-propagation procedure to calculate the directional derivatives w.r.t. input/output values, as described in Sec. 3.4. This results in more memory usage than Adam. The extra computation cost in the CG stage applies to the CG batch, which is only a small portion of the training set. Note that the number of CG iterations used in NG, HF, and NGHF differs in our experiments (see Sec. 7). An extra validation stage is performed after every CG iteration, which increases the computation cost of the proposed methods although it can be performed less frequently with a smaller amount of data. Despite the increased cost, the computation within each gradient batch or CG batch can be easily parallelised (as shown in Fig. 1) which makes it much easier to use many workers with no further approximations than for SGD or Adam. When applied to a much larger training set, it may be unnecessary to scale up the size of the CG batch, which would reduce the relative increase in computation cost.

## 10. Conclusions

A CG-based synchronous distributed optimisation framework for discriminative sequence training has been proposed in this paper, which has the flexibility to combine NG with a second-order optimisation method, such as HF. This framework has the same advantages as HF in providing stable loss value reductions and inherently suitable for parallel computing, yet also has improved numerical issues in CG and improved performance for models with shared parameters. The framework was evaluated using ASR experiments with the training and test data from the MGB challenge. It was shown that the improved CG method can find effective parameter updates resulting in a better improvement in MPE loss often with far fewer CG iterations. Furthermore, when applied in a setting where NG is combined with HF, the resulting NGHF method generates models with better generalisation ability with far fewer parameter updates when compared with SGD and Adam. The method also can be efficiently parallelised across multiple GPUs without making any approximations.

## Acknowledgement

## References

Amari, S., 1997. Neural learning in structured parameter spaces - natural Riemannian gradient, in: Advances in Neural Information Processing Systems 9 (NIPS), pp. 127–133.

Amari, S., 2016. Information Geometry and its Applications. Springer.

Andrychowicz, M., Denil, M., Colmenarejo, S., Hoffman, M., Pfau, D., Schaul, T., Shillingford, B., de Freitas, N., 2016. Learning to learn by gradient descent by gradient descent, in: Advances in Neural Information Processing Systems 29 (NIPS), pp. 3988–3996.

Bahl, L., Brown, P., de Souza, P., Mercer, R., 1986. Maximum mutual information estimation of hidden Markov model parameters for speech recognition, in: Proceedings of the 11th IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 231–234.

Baum, L., Eagon, J., 1967. An inequality with applications to statistical estimation for probabilistic functions of Markov processes and to a model for ecology. Bulletin of the American Mathematical Society 73(3), 360–363.

Becker, S., LeCun, Y., 1988. Improving the Convergence of Back-Propagation learning with Second Order methods. Technical Report CRG-TR-88-5. Computer Science Department, University of Toronto.

Bell, P., Gales, M., Hain, T., Kilgour, J., Lanchantin, P., Liu, X., McParland, A., Renals, S., Saz, O., Wester, M., Woodland, P., 2015. The MGB challenge: Evaluating multi-genre broadcast media recognition, in: Proceedings of the 10th IEEE Workshop on Automatic Speech Recognition and Understanding (ASRU), pp. 687–693.

Bernacchia, A., Lengyel, M., Hennequin, G., 2018. Exact natural gradient in deep linear networks and its application to the nonlinear case, in: Advances in Neural Information Processing Systems 31 (NIPS), pp. 5945–5954.

Bishop, C., 2006. Pattern Recognition and Machine Learning. Springer.

Bottou, L., 2010. Large-scale machine learning with stochastic ´ gradient descent, in: Proceedings of the 19th International Conference on Computational Statistics (COMPSTAT), pp. 177–187.

Bottou, L., Curtis, F.E., Nocedal, J., 2018. Optimization methods for large-scale machine learning. Society for Industrial and Applied Mathematics (SIAM) Review 60(2), 223–311.

Chen, K., Huo, Q., 2016. Scalable training of deep learning machines by incremental block training with intra-block parallel optimization and blockwise model-update filtering, in: Proceedings of the 41st IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 5880–5884.

Chiu, C.C., Sainath, T.N., Wu, Y., Prabhavalkar, R., Nguyen, P., Chen, Z., Kannan, A., Weiss, R.J., Rao, K., Gonina, E., Jaitly, N., Li, B., Chorowski, J., Bacchiani, M., 2018. State-of-the-art speech recognition with sequence-to-sequence models, in: Proceedings of the 43rd IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 4774–4778.

Dean, J., Corrado, G., Monga, R., Chen, K., Devin, M., Le, Q., Mao, M., Ranzato, M., Senior, A., Tucker, P., Yang, K., Ng, A., 2012. Large scale distributed deep networks, in: Advances in Neural Information Processing Systems 25 (NIPS), pp. 1223–1231.

Desjardins, G., Simonyan, K., Pascanu, R., Kavukcuoglu, K., 2015. Natural neural networks, in: Advances in Neural Information Processing Systems 28 (NIPS), pp. 2071–2079.

Dognin, P., Goel, V., 2013. Combining stochastic average gradient and Hessian-free optimization for sequence training of deep neural networks, in: Proc. IEEE Workshop on Automatic Speech Recognition and Understanding (ASRU), pp. 303–308.

Duchi, J., Hazan, E., Singer, Y., 2011. Adaptive subgradient methods for online learning and stochastic optimization. Journal of Machine Learning Research , 2121–2159.

Elman, J., 1990. Finding structure in time. Cognitive Science 14(2), 179–211.

de Felice, F., Clarke, C., 1992. Relativity on Curved Manifolds. Cambridge University Press.

George, T., Laurent, C., Bouthillier, X., Ballas, N., Vincent, P., 2018. Fast approximate natural gradient descent in a Kronecker factored eigenbasis, in: Advances in Neural Information Processing Systems 31 (NIPS), pp. 9550–9560.

Gibson, M., Hain, T., 2006. Hypothesis spaces for minimum Bayes risk training in large vocabulary speech recognition, in: Proceedings of the 7th Conference of the International Speech Communication Association (Interspeech).

Glorot, X., Bengio, Y., 2010. Understanding the difficulty of training deep feedforward neural networks, in: Proceedings of the 13th International Conference on Artificial Intelligence and Statistics (AISTATS), pp. 249–256.

Goel, V., Byrne, W., 2000. Minimum Bayes risk automatic speech recognition. Computer Speech and Language 14(2), 115–135.

Graves, A., Fernández, S., Gomez, F., Schmidhuber, J., 2006. Connectionist temporal classification: Labelling unsegmented sequence data with recurrent neural networks, in: Proceedings of the 23rd International Conference on Machine Learning (ICML), pp. 369–376.

Graves, A., Jaitly, N., 2014. Towards end-to-end speech recognition with recurrent neural networks, in: Proceedings of the 31st International Conference on Machine Learning (ICML), pp. 1764–1772.

Graves, A., Mohamed, A.r., Hinton, G., 2013. Speech recognition with deep recurrent neural networks, in: Proceedings of the 38th IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 6645–6649.

Grosse, R., Salakhudinov, R., 2015. Scaling up natural gradient by sparsely factorizing the inverse Fisher matrix, in: Proceedings of the 32nd International Conference on Machine Learning (ICML), pp. 2304–2313.

Haider, A., 2018. A common framework for natural gradient and Taylor based optimisation using manifold theory. arXiv preprint arXiv:1803.09791.

Haider, A., 2019. Optimisation Methods for Training Deep Neural Networks in Speech Recognition. Ph.D. thesis. University of Cambridge.

Haider, A., Woodland, P., 2017. Sequence training of DNN acoustic models with natural gradient, in: Proceedings of the 11th IEEE Workshop on Automatic Speech Recognition and Understanding (ASRU), pp. 178–184.

Haider, A., Woodland, P., 2018. Combining natural gradient with Hessian free methods for sequence training, in: Proceedings of the 19th Conference of the International Speech Communication Association, pp. 2918–2922.

Heigold, G., McDermott, E., Vanhoucke, V., Senior, A., Bacchiani, M., 2014. Asynchronous stochastic optimization for sequence training of deep neural networks, in: Proceedings of the 39th IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 5624–5628.

Hinton, G., Deng, L., Yu, D., Dahl, G.E., Mohamed, A.r., Jaitly, N., Senior, A., Vanhoucke, V., Nguyen, P., Sainath, T., Brian, K., 2012. Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups. IEEE Signal Processing Magazine 29(6), 82–97.

Hinton, G., Osindero, S., Teh, Y., 2006. A fast learning algorithm for deep belief nets. Neural Computation 18, 1527–1554.

Hinton, G., Salakhutdinov, R., 2006. Reducing the dimensionality of data with neural networks. Science 313(5786), 504–507.

Hochreiter, S., Schmidhuber, J., 1997. Long short-term memory. Neural computation 9(8), 1735–1780.

Kaiser, J., Horvat, B., Kacic, Z., 2000. A novel loss function for the overall risk criterion based discriminative training of HMM models, in: Proceedings of the 1st Conference of the International Speech Communication Association (Interspeech), pp. 887–890.

Kingma, D.P., Ba, J.L., 2015. Adam: A method for stochastic optimization, in: Proceedings of the 3rd International Conference on Learning Representations (ICLR), pp. 1–13.

Kingsbury, B., 2009. Lattice-based optimization of sequence classification criteria for neural-network acoustic modeling, in: Proceedings of the 34th IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 3761–3764.

Kingsbury, B., Sainath, T., Soltau, H., 2012. Scalable minimum Bayes risk training of deep neural networks acoustic models using distributed Hessian-free optimization, in: Proceedings of the 13th Conference of the International Speech Communication Association (Interspeech), pp. 10–13.

Kreyssig, F., Zhang, C., Woodland, P., 2018. Improved TDNNs using Deep Kernels and Frequency Dependent Grid-RNNs, in: 2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 4864–4868.

Ladkat, P., Rybakov, O., Arava, R., Hari, S., Parthasarathi, K., Chen, I.F., Ström, N., 2019. Two tiered distributed training algorithm for acoustic modeling, in: Proceedings of the 20th Conference of the International Speech Communication Association (Interspeech), pp. 1626–1630.

Lüscher, C., Beck, E., Irie, K., Kitza, M., Michel, W., Zeyer, A., Schlüter, R., Ney, H., 2019. RWTH ASR systems for LibriSpeech: Hybrid vs attention, in: Proceedings of the 20th Conference of the International Speech Communication Association (Interspeech), pp. 231–235.

Martens, J., 2010. Deep learning via Hessian-free optimization, in: Proceedings of the 27th International Conference on Machine Learning (ICML), pp. 735–742.

Martens, J., 2020. New insights and perspectives on the natural gradient method. Journal of Machine Learning Research , 1–76.

Martens, J., Grosse, R., 2015. Optimizing neural networks with Kronecker-factored approximate curvature, in: Proceedings of the 32nd International Conference on Machine Learning (ICML), pp. 2408–2417.

Martens, J., Sutskever, I., 2011. Learning recurrent neural networks with Hessian-free optimization, in: Proceedings of the 28th International Conference on Machine Learning (ICML), pp. 1033–1040.

Nocedal, J., Wright, S., 2016. Numerical Optimization. Springer.

Pascanu, R., Bengio, Y., 2013. Revisiting natural gradient for deep networks. arXiv preprint arXiv:1301.3584.

Pascanu, R., Mikolov, T., Bengio, Y., 2013. On the difficulty of training recurrent neural networks, in: Proceedings of the 30th International Conference on Machine Learning (ICML), pp. 1310–1318.

Pearlmutter, B., 1994. Fast Exact Multiplication by the Hessian. Neural computation 6(1), 147–160.

Peddinti, V., Povey, D., Khudanpur, S., 2015. A time delay neural network architecture for efficient modeling of long temporal contexts, in: Proceedings of the 16th Conference of the International Speech Communication Association (Interspeech), pp. 3214–3218.

Povey, D., 2005. Discriminative training for Large Vocabulary Speech Recognition. Ph.D. thesis. University of Cambridge.

Povey, D., Peddinti, V., Galvez, D., Ghahremani, P., Manohar, V., Na, X., Wang, Y., Khudanpur, S., 2016. Purely sequence-trained neural networks for ASR based on lattice-free MMI, in: Proceedings of the 17th Conference of the International Speech Communication Association (Interspeech).

Povey, D., Woodland, P., 2002. Minimum phone error and I-smoothing for improved discriminative training, in: Proceedings of the 27th IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 105–108.

Povey, D., Zhang, X., Khudanpur, S., 2015. Parallel training of deep neural networks with natural gradient and parameter averaging. Proceedings of the 3rd International Conference on Learning Representations (ICLR) , 1–13.

Renals, S., Morgan, N., Cohen, M., Franco, H., 1992. Connectionist probability estimation in the DECIPHER speech recognition system, in: Proceedings of the 17th IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 601–604.

Robinson, A., Fallside, F., 1987. The utility driven dynamic error propagation network. Technical Report CUED/F-INFENG/TR.1. Cambridge University Engineering Department.

Roux, N., Manzagol, P.A., Bengio, Y., 2008. Topmoumoute online natural gradient algorithm, in: Advances in Neural Information Processing Systems 20 (NIPS), pp. 849–856.

Rumelhart, D., Hinton, G., Williams, R., 1986. Learning representations by back-propagating errors. Nature 323(6088), 533–536.

Sainath, T., Horesh, L., Kingsbury, B., Aravkin, A., Ramabhadran, B., 2013a. Accelerating Hessian-free optimization for deep neural networks by implicit preconditioning and sampling, in: Proceedings of the 9th IEEE Workshop on Automatic Speech Recognition and Understanding (ASRU), pp. 303–308.

Sainath, T., Kingsbury, B., Soltau, H., Ramabhadran, B., 2013b. Optimization techniques to improve training speed of deep neural networks for large speech tasks. IEEE Transactions on Audio, Speech, and Language Processing 21(11), 2267–2276.

Sak, H., Senior, A., Beaufays, F., 2014. Long short-term memory recurrent neural network architectures for large scale acoustic modeling, in: Proceedings of the 15th Conference of the International Speech Communication Association (Interspeech), pp. 338–342.

Saon, G., Sercu, T., Rennie, S., Kuo, H.K., 2016. The IBM 2016 English

conversational telephone speech recognition system, in: Proceedings of the 17th Conference of the International Speech Communication Association (Interspeech), pp. 7–11.

Schraudolph, N., 2002. Fast curvature matrix-vector products for second-order gradient descent. Neural Computation 14(7), 1723–1738.

Schulman, J., Levine, S., Moritz, P., Jordan, M., Abbeel, P., 2015. Trust region policy optimization, in: Proceedings of the 31st International Conference on Machine Learning (ICML), pp. 889–1897.

Seide, F., Fu, H., Droppo, J., Li, G., Yu, D., 2014. 1-bit stochastic gradient descent and its application to data-parallel distributed training of speech DNNs, in: Proceedings of the 15th Conference of the International Speech Communication Association, pp. 1058–1062.

Senior, A., Heigold, G., Ranzato, M., Yang, K., 2013. An empirical study of learning rates in deep neural networks for speech recognition, in: Proceedings of the 38th IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 6724–6728.

Shannon, M., 2017. Optimizing expected word error rate via sampling for speech recognition. Proceedings of the 18th Conference of the International Speech Communication Association (Interspeech) .

Shewchuk, J., 1994. An Introduction to the Conjugate Gradient Method without the Agonizing Pain. Technical Report CMU-CS-94-125. Department of Computer Science,Carnegie-Mellon University.

Simonyan, K., Zisserman, A., 2015. Very deep convolutional networks for large-scale image recognition, in: Proceedings of the 3rd International Conference on Learning Representations (ICLR), pp. 1–14.

Ström, N., 2015. Scalable distributed DNN training using commodity GPU cloud computing, in: Proceedings of the 16th Conference of the International Speech Communication Association (Interspeech), pp. 1488–1492.

Su, H., Li, G., Yu, D., Seide, F., 2013. Error back propagation for sequence training of context-dependent deep neural networks for conversational speech transcription, in: Proceedings of the 38th IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 6664–6668.

Tieleman, T., Hinton, G., 2012. Lecture 6.5–RMSprop: Divide the gradient by a running average of its recent magnitude. COURSERA: Neural Networks for Machine Learning.

Tikhonov, A., Leonov, A., Yagola, A., 1998. Nonlinear Ill-Posed Problems. Springer.

Valtchev, V., 1995. Discriminative Methods in HMM-based Speech Recognition. Ph.D. thesis. University of Cambridge.

Valtchev, V., Odell, J., Woodland, P., Young, S., 1997. MMIE training of large vocabulary recognition systems. Speech Communication 22(4), 303–314.

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A., Kaiser, L., I., P., 2017. Attention is all you need, in: Advances in Neural Information Processing Systems 30 (NIPS), pp. 6000–6010.

Veselỳ, K., Ghoshal, A., Burget, L., Povey, D., 2013. Sequence-discriminative training of deep neural networks, in: Proceedings of the 14th Conference of the International Speech Communication Association (Interspeech).

Vinyals, O., Povey, D., 2012. Krylov subspace descent for deep learning, in: Proceedings of the 15th International Conference on Artificial Intelligence and Statistics (AISTATS), pp. 1261–1268.

Waibel, A., Hanazawa, T., Hinton, G., Shikano, K., Lang, J., 1989. Phoneme recognition using time-delay neural networks. IEEE Transactions on Acoustic, Speech, and Signal Processing 37(3), 328–339.

Werbos, P., 1988. Generalization of backpropagation with application to a recurrent gas market model. Neural computation 1(4), 339–356.

Wiesler, S., Golik, P., Schlüter, R., Ney, H., 2015. Investigations on sequence training of neural networks, in: Proceedings of the 40th IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 4565–4569.

Wiesler, S., Li, J., Xue, J., 2013. Investigations on Hessian-free optimization for cross-entropy training of deep neural networks, in: Proceedings of the 14th Conference of the International Speech Communication Association (Interspeech), pp. 3317–3321.

Woodland, P., Liu, X., Qian, Y., Zhang, C., Gales, M., Karanasou, P., Lan-

chantin, P., Wang, L., 2015. Cambridge University transcription systems for the Multi-Genre Broadcast challenge, in: Proceedings of the 10th IEEE Workshop on Automatic Speech Recognition and Understanding (ASRU), pp. 639–646.

Woodland, P., Povey, D., 2002. Large scale discriminative training of hidden Markov models for speech recognition. Computer Speech and Language 16(1), 25–47.

Xiong, W., Droppo, J., Huang, X., Seide, F., Seltzer, M., Stolcke, A., Yu, D., Zweig, G., 2016. The microsoft 2016 conversational speech recognition system, in: Proceedings of the 41st IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 5255–5259.

Xu, P., Roosta, F., Mahoney, M., 2020. Second-order optimization for nonconvex machine learning: An empirical study, in: Proceedings of the 2020 SIAM International Conference on Data Mining (SDM), pp. 199–207.

Young, S., Evermann, G., Gales, M., Hain, T., Kershaw, D., Liu, X., Moore, G., Odell, J., Ollason, D., Povey, D., Ragni, A., Valtchev, V., Woodland, P., Zhang, C., 2015. The HTK Book (for HTK version 3.5). Cambridge University Engineering Department.

Zeiler, M., 2012. ADADELTA: An adaptive learning rate method. arXiv preprint `arXiv:1212.5701`.

Zhang, C., 2017. Joint Training Methods for Tandem and Hybrid Speech Recognition Systems using Deep Neural Networks. Ph.D. thesis. University of Cambridge.

Zhang, C., Kreyssig, F., Li, Q., Woodland, P., 2019a. PyHTK: Python library and ASR pipelines for HTK, in: Proc. ICASSP, pp. 6470–6474.

Zhang, C., Woodland, P., 2015. A general artificial neural network extension for HTK, in: Proceedings of the 16th Conference of the International Speech Communication Association, pp. 3581–3585.

Zhang, C., Woodland, P., 2017. Joint optimisation of tandem systems using Gaussian mixture density neural network discriminative sequence training, in: Proceedings of the 42nd IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 5015–5019.

Zhang, W., Cui, X., Finkler, U., Kingsbury, B., Saon, G., Kung, D., Picheny, M., 2019b. Distributed deep learning strategies for automatic speech recognition, in: Proceedings of the 44th IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 5706–5710.

Zhang, W., Cui, X., Finkler, U., Saon, G., Kayi, A., Buyuktosunoglu, A., Kingsbury, B., Kung, D., Picheny, M., 2019c. A highly efficient distributed deep learning system for automatic speech recognition, in: Proceedings of the 20th Conference of the International Speech Communication Association (Interspeech), pp. 2628–2632.

Zinkevich, M., Weimer, M., Li, L., Smola, A., 2010. Parallelized stochastic gradient descent, in: Advances in Neural Information Processing Systems 23 (NIPS), pp. 2595–2603.