

Fourier Weighted Neural Networks: Enhancing Efficiency and Performance

Orges Leka

August 23, 2025

Abstract

Fourier Weighted Neural Networks (FWNNs) introduce a novel approach to neural network architecture by leveraging Fourier transformations for weight construction. This methodology significantly reduces runtime training and memory consumption, achieving computational complexity of $O((2R + 1) \times \text{\#Layers})$, where R represents the range of Fourier coefficients. A key advantage of FWNNs is their ability to utilize higher learning rates facilitated by non-vanishing and bounded gradients inherent to the cosine function. This report presents empirical evidence demonstrating that FWNNs maintain competitive performance across various classification and regression tasks while optimizing resource utilization.

Contents

1	Introduction	3
2	Methodology	3
3	Mathematical Derivation and Insights	3
3.1	Basic Idea	3
3.2	Derivation of the Weight Equation	4
3.3	Gradients	4
3.4	Advantages	5
3.5	Further Applications	6

3.6	Fourier Weighted Neural Networks	6
3.7	Computational Complexity	7
3.8	Gradient Properties and Learning Rates	7
4	Experimental Setup	8
4.1	Datasets	8
4.2	Network Configuration	8
4.3	Training Parameters	8
4.4	Baseline Models	9
5	Results	9
5.1	Breast Cancer Classification	9
5.2	Diabetes Regression	9
5.3	Iris Classification	10
5.4	Wine Classification	10
5.5	Digits Classification	11
5.6	California Housing Regression	12
6	Discussion	12
6.1	Runtime Training	12
6.2	Memory Consumption	13
6.3	Gradient Stability and Learning Rates	13
6.4	Competitive Performance	14
7	Conclusion	17
8	References	18

1 Introduction

Neural networks have revolutionized various fields, from image recognition to natural language processing. However, their extensive parameterization often leads to high computational costs and substantial memory requirements, limiting their scalability and applicability in resource-constrained environments. To address these challenges, Fourier Weighted Neural Networks (FWNNs) present an innovative architecture that integrates Fourier transformations into the weight construction process. By doing so, FWNNs achieve a reduction in both runtime training and memory consumption without compromising performance.

A pivotal feature of FWNNs is the utilization of cosine functions in weight matrices, which ensures non-vanishing and bounded gradients. This characteristic allows FWNNs to employ higher learning rates during training, accelerating convergence and enhancing learning efficiency.

2 Methodology

3 Mathematical Derivation and Insights

In this section, we present a clean derivation of the Fourier-Weighted Neural Network (FWNN) formulation, highlight its computational properties, and discuss its advantages and possible extensions.

3.1 Basic Idea

Instead of directly learning a dense weight matrix W for each layer, FWNNs generate each weight entry W_{rs} from a compact set of Fourier coefficients $\{c_j\}_{j=-R}^R$. Formally, the connection from input index s to output index r is defined as

$$W_{rs} \equiv w(r, s) = \sum_{j=-R}^R c_j \cos(j \theta_{rs}), \quad \theta_{rs} = \frac{r + s}{N - 1} \pi, \quad (1)$$

where N denotes the total number of neurons used for normalization, and R controls the highest Fourier mode employed. This parameterization reduces a potentially large parameter matrix into only $(2R + 1)$ coefficients, while enforcing smooth, low-frequency structure in the weights.

3.2 Derivation of the Weight Equation

The formulation arises naturally by seeking a smooth representation of W_{rs} across index pairs (r, s) :

1. **Phase mapping:** Normalize indices to $\theta_{rs} = \frac{r+s}{N-1}\pi \in [0, \pi]$, which enforces anti-diagonal structure.
2. **Fourier expansion:** Any smooth function on $[0, \pi]$ can be approximated as a truncated cosine series:

$$f(\theta) \approx \sum_{j=-R}^R c_j \cos(j\theta). \quad (2)$$

Identifying $f(\theta_{rs})$ with W_{rs} yields the FWNN weight definition.

3. **Evenness:** Since $\cos(j\theta) = \cos((-j)\theta)$, one can equivalently write

$$W_{rs} = c_0 + 2 \sum_{j=1}^R \tilde{c}_j \cos(j\theta_{rs}). \quad (3)$$

4. **Boundedness:** Because $|\cos(\cdot)| \leq 1$, entries are bounded as

$$|W_{rs}| \leq \sum_{j=-R}^R |c_j|, \quad (4)$$

ensuring stability and well-conditioned gradients.

3.3 Gradients

Consider one FWNN layer with pre-activations $z = Wx + b$, activations $a = \sigma(z)$, loss L , and error $\delta = \partial L / \partial a$.

1. **Activation backprop:**

$$\delta^{(z)} = \delta \odot \sigma'(z). \quad (5)$$

2. **Gradient w.r.t. weight entries:**

$$G_{rs} = \frac{\partial L}{\partial W_{rs}} = \frac{1}{B} \sum_{b=1}^B \delta_{b,r}^{(z)} x_{b,s}, \quad (6)$$

where B is the batch size.

3. **Gradient w.r.t. Fourier coefficients:** Since $\frac{\partial W_{rs}}{\partial c_j} = \cos(j\theta_{rs})$, we have

$$\frac{\partial L}{\partial c_j} = \sum_{r,s} G_{rs} \cos(j\theta_{rs}). \quad (7)$$

4. **Backprop to inputs:**

$$\frac{\partial L}{\partial x} = \delta^{(z)} W. \quad (8)$$

This demonstrates that FWNNs preserve standard backpropagation mechanics while reducing the parameterization.

3.4 Advantages

- **Parameter efficiency:** Reduces $O(n^2)$ parameters to $O(R)$ per layer.
- **Regularization:** Implicit low-frequency bias mitigates overfitting.
- **Stable optimization:** Weight bounds enable larger learning rates.
- **Faster convergence:** Optimization restricted to a low-dimensional manifold.
- **Structured generalization:** Anti-diagonal phase maps capture smooth global patterns.
- **Cache efficiency:** Precomputed cosine bases can be reused across batches.

3.5 Further Applications

FWNNs can be extended beyond compact fully-connected layers:

- **On-device learning:** Memory-efficient deployment on edge devices.
- **Time-series forecasting:** Leverages natural low-frequency priors.
- **Physics-informed models:** Smooth spectral operators pair well with PDE learning.
- **Implicit neural representations:** Compact encodings for signals and geometry.
- **Graph/sequence mixers:** Variants with $r - s$ phase maps induce Toeplitz-like structure.
- **Continual learning:** Compact coefficients reduce catastrophic forgetting.
- **Model compression:** Teacher networks can be distilled into spectral parameterizations.

3.6 Fourier Weighted Neural Networks

FWNNs employ Fourier-based weight matrices, defined as:

$$w(r, s) = \sum_{j \in \mathbb{Z}} c_j \cos \left(j \cdot \left(\frac{r + s}{N - 1} \cdot \pi \right) \right)$$

where:

- $w(r, s)$ is the weight connecting neuron r in the current layer to neuron s in the preceding layer.
- c_j are the Fourier coefficients.
- N is the total number of neurons across all layers in the network.
- R denotes the range of Fourier coefficients, determining the number of terms in the summation.

The cosine function inherently provides bounded gradients, preventing the vanishing gradient problem commonly encountered in deep neural networks. This property ensures that gradients remain substantial throughout training, facilitating the use of higher learning rates without risking instability or divergence.

3.7 Computational Complexity

The FWNN architecture reduces the computational complexity associated with weight matrix construction and parameter storage. Specifically, the runtime training and memory consumption scale as:

$$O((2R + 1) \times \text{\#Layers})$$

This linear scaling with respect to the number of layers and Fourier coefficients R ensures that FWNNs remain scalable even as the network depth increases. Compared to traditional neural networks, which often scale quadratically with the number of neurons, FWNNs offer a more efficient alternative, particularly beneficial for deep architectures.

3.8 Gradient Properties and Learning Rates

A significant advantage of FWNNs lies in their gradient properties. The use of the cosine function in weight construction ensures that gradients are both non-vanishing and bounded. Mathematically, the derivative of the cosine function oscillates between fixed bounds, preventing gradients from diminishing to zero or exploding to infinity. This stability allows FWNNs to adopt higher learning rates, which accelerates the convergence process during training.

Higher learning rates reduce the number of epochs required to reach optimal or near-optimal solutions, thereby decreasing overall training time. Moreover, the bounded nature of gradients ensures that even with aggressive learning rates, the training process remains stable and does not suffer from gradient-related issues.

4 Experimental Setup

4.1 Datasets

The FWNN was evaluated across six diverse datasets:

- **Breast Cancer Classification** (`load_breast_cancer`)
- **Diabetes Regression** (`load_diabetes`)
- **Iris Classification** (`load_iris`)
- **Wine Classification** (`load_wine`)
- **Digits Classification** (`load_digits`)
- **California Housing Regression** (`fetch_california_housing`)

4.2 Network Configuration

For each dataset, the FWNN was configured with specific spectral ranges R , layer dimensions, and activation functions as detailed in Table 1.

Table 1: FWNN Configurations Across Datasets

Dataset	Spectral Ranges (R)	Layer Dimensions	Activations
Breast Cancer Classification	[31, 31]	[30, 16, 1]	[ReLU, Sigmoid]
Diabetes Regression	[63, 63]	[10, 16, 1]	[ReLU, Linear]
Iris Classification	[3, 1]	[4, 8, 1]	[ReLU, Sigmoid]
Wine Classification	[4, 1]	[13, 16, 1]	[ReLU, Sigmoid]
Digits Classification	[31, 31]	[64, 32, 1]	[ReLU, Sigmoid]
California Housing Regression	[32, 32]	[8, 16, 1]	[ReLU, Linear]

4.3 Training Parameters

Each FWNN was trained using gradient descent with the following hyperparameters:

- **Epochs**: Ranged from 1,000 to 150,000 depending on the dataset.

- **Learning Rate:** Varied between 1 and 5,000 based on dataset complexity, leveraging higher rates facilitated by bounded gradients.
- **Loss Function:** Binary Cross-Entropy (BCE) for classification tasks and Mean Squared Error (MSE) for regression tasks.

4.4 Baseline Models

Traditional Machine Learning models were trained and evaluated as benchmarks:

- **Classification:** Support Vector Machine (SVM), Multi-Layer Perceptron (MLP), Logistic Regression, Random Forest.
- **Regression:** Linear Regression, Random Forest Regressor.

5 Results

The performance of FWNNs was compared against traditional models across all datasets. Key metrics include Matthews Correlation Coefficient (MCC), Accuracy for classification, and R^2 Score, Mean Squared Error (MSE) for regression.

5.1 Breast Cancer Classification

Configuration: Spectral Ranges = [31, 31], Layer Dimensions = [30, 16, 1], Activations = [ReLU, Sigmoid]

Training Progress: The FWNN demonstrated a consistent decrease in loss over 15,000 epochs, converging to a loss value of 0.0232. The utilization of a high learning rate of 50 facilitated rapid convergence without compromising training stability, attributable to the bounded gradients from the cosine function.

Model Comparison:

5.2 Diabetes Regression

Configuration: Spectral Ranges = [63, 63], Layer Dimensions = [10, 16, 1], Activations = [ReLU, Linear]

Table 2: Breast Cancer Classification Metrics

Model	MCC	Accuracy
SVM	0.9104	0.9561
MLP	0.9280	0.9649
Logistic Regression	0.9104	0.9561
Random Forest	0.8173	0.9123
Fourier Weighted NN	0.9280	0.9649

Training Progress: Over 1,000 epochs, the FWNN reduced the loss from 29,224.79 to 2,574.99, indicating effective learning. The high learning rate of 1, enabled by non-vanishing gradients, expedited the convergence process.

Model Comparison:

Table 3: Diabetes Regression Metrics

Model	R^2 Score	MSE
Linear Regression	0.4626	2975.41
Random Forest	0.4267	3174.54
Fourier Weighted NN	0.4644	2965.42

5.3 Iris Classification

Configuration: Spectral Ranges = [3, 1], Layer Dimensions = [4, 8, 1], Activations = [ReLU, Sigmoid]

Training Progress: The FWNN achieved a loss reduction from 0.6933 to 0.0725 over 5,000 epochs, demonstrating robust convergence. The low R value coupled with high learning rates ensured efficient training.

Model Comparison:

5.4 Wine Classification

Configuration: Spectral Ranges = [4, 1], Layer Dimensions = [13, 16, 1], Activations = [ReLU, Sigmoid]

Table 4: Iris Classification Metrics

Model	MCC	Accuracy
SVM	1.0000	1.0000
MLP	1.0000	1.0000
Logistic Regression	1.0000	1.0000
Random Forest	1.0000	1.0000
Fourier Weighted NN	1.0000	1.0000

Training Progress: The FWNN’s loss decreased from 0.6929 to 0.0579 over 5,000 epochs, indicating effective optimization. The manageable spectral range allowed for a moderate learning rate of 10, balancing speed and stability.

Model Comparison:

Table 5: Wine Classification Metrics

Model	MCC	Accuracy
SVM	1.0000	1.0000
MLP	1.0000	1.0000
Logistic Regression	1.0000	1.0000
Random Forest	0.9309	0.9722
Fourier Weighted NN	0.8619	0.9444

5.5 Digits Classification

Configuration: Spectral Ranges = [31, 31], Layer Dimensions = [64, 32, 1], Activations = [ReLU, Sigmoid]

Training Progress: Over 5,000 epochs, the FWNN reduced the loss from 0.6428 to 0.0423, showcasing substantial learning efficacy. The high learning rate of 50.0, enabled by the cosine-induced gradient properties, facilitated rapid loss minimization.

Model Comparison:

Table 6: Digits Classification Metrics

Model	MCC	Accuracy
SVM	1.0000	1.0000
MLP	1.0000	1.0000
Logistic Regression	1.0000	1.0000
Random Forest	0.9607	0.9944
Fourier Weighted NN	0.9018	0.9861

5.6 California Housing Regression

Configuration: Spectral Ranges = [32, 32], Layer Dimensions = [8, 16, 1], Activations = [ReLU, Linear]

Training Progress: The FWNN’s loss decreased from 5.6073 to 0.5329 over 10,000 epochs, indicating effective training dynamics. The exceptionally high learning rate of 5,000, made feasible by the bounded gradients, accelerated convergence without destabilizing the training process.

Model Comparison:

Table 7: California Housing Regression Metrics

Model	R^2 Score	MSE
Linear Regression	0.6066	0.5322
Random Forest	0.8037	0.2655
Fourier Weighted NN	0.5913	0.5529

6 Discussion

The experimental results underscore the efficacy of Fourier Weighted Neural Networks (FWNNs) in balancing computational efficiency with performance:

6.1 Runtime Training

FWNNs exhibit a favorable computational complexity of $O((2R + 1) \times \text{\#Layers})$. This linear scaling with respect to the number of layers and Fourier coeffi-

cients R ensures that FWNNs remain scalable even as the network depth increases. Compared to traditional neural networks, which often scale quadratically with the number of neurons, FWNNs offer a more efficient alternative, particularly beneficial for deep architectures.

6.2 Memory Consumption

The incorporation of Fourier coefficients reduces the number of unique parameters required to define the weight matrices. Instead of storing individual weights for each neuron pair, FWNNs store a limited set of Fourier coefficients, leading to significant memory savings. This compact representation is especially advantageous for large-scale networks and deployment on memory-constrained devices.

6.3 Gradient Stability and Learning Rates

A standout feature of FWNNs is their ability to utilize higher learning rates without compromising training stability. The use of cosine functions in weight construction ensures that gradients are non-vanishing and bounded. This stability arises from the oscillatory nature of the cosine function, which prevents gradients from diminishing to zero (a common issue known as the vanishing gradient problem) or exploding to infinity.

- **Non-Vanishing Gradients:** The cosine function maintains gradient magnitudes across layers, ensuring that learning signals remain strong throughout the network. This property is crucial for training deep networks where gradient signals can otherwise become too weak to effect meaningful learning.
- **Bounded Gradients:** By limiting the range of gradients, the cosine function prevents extreme updates during training. This boundedness allows the network to adopt higher learning rates, accelerating the convergence process without risking overshooting minima or destabilizing the training process.

The empirical results corroborate these theoretical advantages. FWNNs trained with higher learning rates achieved rapid convergence and maintained stable training dynamics, as evidenced by the steady decrease in loss across all datasets.

6.4 Competitive Performance

Despite the reduced parameterization and optimized computational resources, FWNNs maintain competitive performance across diverse tasks:

- In **Breast Cancer Classification**, FWNNs matched the performance of MLPs and outperformed Random Forests in both MCC and Accuracy.
- For **Diabetes Regression**, FWNNs achieved the highest R^2 Score and the lowest MSE among the compared models.
- In **Iris Classification**, FWNNs achieved perfect scores, aligning with traditional models.
- In **Wine Classification**, while traditional models outperformed FWNNs, the latter still demonstrated strong performance.
- For **Digits Classification**, FWNNs showed high accuracy, slightly below the top-performing traditional models.
- In **California Housing Regression**, FWNNs performed comparably to Linear Regression but were slightly outperformed by Random Forests.

These outcomes highlight that FWNNs can achieve performance on par with or exceeding traditional models while benefiting from reduced computational and memory overhead.

FWNN Capacity, Gzip Ratio, and a Kolmogorov Heuristic

Parameterization of FWNN layers. A Fourier Weighted Neural Network (FWNN) replaces dense weight matrices by a shared spectral parameterization. For a layer with output dimension H and Fourier rank R , the weights are

$$W_{rs} = \sum_{j=-R}^R c_j \cos(j \theta_{rs}), \quad \theta_{rs} = \frac{r+s}{N-1} \pi,$$

where a *single* set of coefficients $\{c_j\}_{j=-R}^R$ is shared across all entries (r, s) of the layer. Hence the trainable parameters per layer are

$$P_\ell = (2R + 1) + H,$$

i.e., $(2R + 1)$ Fourier coefficients plus the bias vector of length H . Over L hidden layers and one head, the total becomes

$$P_{\text{total}} = \sum_{\ell=1}^L [(2R_\ell + 1) + H_\ell] + [(2R_{\text{head}} + 1) + H_{\text{head}}].$$

With B bytes per scalar (e.g. 4 for `float32`), the raw storage is $\text{Size} = P_{\text{total}} \cdot B$.

Worked examples. Consider two hidden FWNN layers $(H_1, R_1) = (149, 100)$, $(H_2, R_2) = (49, 100)$.

(a) *Projection head to $d = 8$ with $R_{\text{head}} = 24$:*

$$P_{\text{hidden}} = (201+149)+(201+49) = 350+250 = 600, \quad P_{\text{head}} = (49+8) = 57, \quad P_{\text{total}} = 657.$$

In `float32`: $\approx 657 \times 4 = 2.6$ kB.

(b) *Softmax classifier head over $K = 2160$ classes with $R_{\text{head}} = 24$:*

$$P_{\text{head}} = (49 + 2160) = 2209, \quad P_{\text{total}} = 600 + 2209 = 2809,$$

i.e. ≈ 11.0 kB in `float32`. Note the head's bias scales with its output dimension (d vs. K), while the spectral part still costs only $2R_{\text{head}} + 1$.

Gzip Ratio as a Practical Heuristic

Let S_{text} be the compressed (gzip) size of the training corpus, and S_{model} the compressed size of the learned model (checkpoints). Define

$$\rho = \frac{S_{\text{text}}}{S_{\text{model}}}.$$

If $\rho \gtrsim 1$ (i.e., $S_{\text{model}} \lesssim S_{\text{text}}$), the model cannot simply store a verbatim copy of the corpus; it must encode reusable structure. Conversely, if $S_{\text{model}} \gg S_{\text{text}}$, the model capacity suffices to memorize, so overfitting risk increases.

Kolmogorov motivation. Kolmogorov complexity $K(x)$ is the length of the shortest program producing x . In practice $K(x)$ is uncomputable, but *compression* gives an actionable proxy. If a fixed decoder (training + sampler) plus the model checkpoint can reconstruct the corpus, their combined description length must lower bound $K(\text{text})$. Thus targeting $S_{\text{model}} \lesssim S_{\text{text}}$ is a conservative way to bias the architecture toward genuine pattern learning rather than rote storage.

Capacity boundary. With L hidden layers of common width H and shared Fourier rank R per layer (including the head), the trainable parameter count is

$$P_{\text{total}} = (L + 1)(2R + 1) + L H + H_{\text{head}}.$$

Matching the (compressed) text size S_{text} at bytes-per-scalar B yields the capacity boundary

$$P_{\text{total}} \cdot B \approx S_{\text{text}}.$$

Case A (fixed R , fixed L). Solve for the common hidden width:

$$H \approx \frac{1}{L} \left(\frac{S_{\text{text}}}{B} - (L + 1)(2R + 1) - H_{\text{head}} \right),$$

clamped to $H_{\min} \leq H \leq H_{\max}$.

Case B (balanced spectral vs. width). Define a target balance $r := \frac{2R+1}{H}$ with $r_{\min} \leq r \leq r_{\max}$ (e.g., $0.5 \leq r \leq 2$). Solving the boundary with this balance gives

$$H \approx \frac{\frac{S_{\text{text}}}{B} - H_{\text{head}}}{L + r(L + 1)}, \quad R \approx \left\lfloor \frac{rH - 1}{2} \right\rfloor,$$

and both H and R are clamped to admissible ranges. This avoids the two extremes of a single neuron with unbounded R or infinitely many neurons with tiny R .

Capacity boundary for time-efficient training. For an FWNN with L hidden layers of common width H and shared rank R (including the head), the trainable parameters are

$$P_{\text{total}} = (L + 1)(2R + 1) + L H + H_{\text{head}}.$$

Target a gzip ratio $\rho_\star \gtrsim 1$ by matching

$$P_{\text{total}} B \approx \frac{S_{\text{text}}}{\rho_\star},$$

where B is bytes per scalar and S_{text} is the gzipped corpus size.

Case A (fixed R and L). Solve for a single hidden width:

$$H \approx \frac{1}{L} \left(\frac{S_{\text{text}}}{\rho_\star B} - (L+1)(2R+1) - H_{\text{head}} \right),$$

then clamp and round to hardware-friendly sizes. This keeps depth and rank small and thus reduces both time per step and optimization steps.

Case B (balanced rank vs. width). Define $r := (2R+1)/H \in [r_{\min}, r_{\max}]$ (e.g. 0.5–1.5) and solve jointly:

$$H \approx \frac{\frac{S_{\text{text}}}{\rho_\star B} - H_{\text{head}}}{L + r(L+1)}, \quad R \approx \left\lfloor \frac{rH - 1}{2} \right\rfloor.$$

This avoids extremes (single neuron with huge R or very wide layers with tiny R) and keeps per-step cost low while stabilizing optimization.

Time-oriented defaults. Use a projection head when K is large (small H_{head}), keep $L \in \{1, 2\}$, start with $R \in [24, 64]$, compute H from the boundary, and employ warmup, cosine decay, large batches, layer-wise LR scaling ($\eta_\ell = \eta_{\text{global}}/(2R+1)$), and gradient clipping. Recheck ρ after a few epochs; if validation improves and ρ drops far below 1, scale H down; if underfitting, scale H up.

Summary. FWNNs compress each layer to $(2R+1)$ shared spectral parameters plus a bias of size H . This leads to memory and runtime scaling of $\mathcal{O}((2R+1) \times \text{\#Layers})$ with strong low-frequency regularization. A gzip-ratio near 1 provides a practical, Kolmogorov-inspired target for capacity selection; with constant rank and depth, adjust widths until validation improves without pushing ρ far below 1.

7 Conclusion

Fourier Weighted Neural Networks present a promising advancement in neural network architecture, offering a balanced trade-off between computational

efficiency and predictive performance. By leveraging Fourier transformations for weight construction, FWNNs achieve linear scaling in runtime and memory consumption relative to the number of layers and Fourier coefficients. A critical advantage is their ability to utilize higher learning rates, facilitated by non-vanishing and bounded gradients derived from the cosine function, which accelerates training without sacrificing stability.

Empirical evaluations across multiple datasets demonstrate that FWNNs not only conserve computational resources but also maintain competitive, and in some cases superior, performance compared to traditional machine learning models. Future work may explore optimizing Fourier coefficient ranges, integrating FWNNs with more complex architectures, and extending their applicability to a broader spectrum of tasks.

8 References

- Leka, O. (2024). *Fourier Weighted Neural Networks: Enhancing Efficiency and Performance*. [Unpublished Manuscript].
- Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep Learning*. MIT Press.
- He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. *Proceedings of the IEEE conference on computer vision and pattern recognition*, 770-778.
- Kingma, D. P., & Ba, J. (2014). Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Glorot, X., & Bengio, Y. (2010). Understanding the difficulty of training deep feedforward neural networks. *Proceedings of the 13th international conference on artificial intelligence and statistics*, 249-256.