# UNIT II

# Difficulty of training DNNs

- Deep Neural Networks are very powerful and can do amazing things, but training them can be difficult.

- Training deep neural networks (DNNs) has led to impressive advances in artificial intelligence.

- However, it comes with hurdles like vanishing gradients, overfitting, and limited labeled data.

- **Vanishing and Exploding Gradients**

- Deep learning networks can be problematic when the numbers change too quickly or slowly through many layers.

- This can make it hard for the network to learn and stay stable.

- This can cause difficulties for the network in learning and remaining stable.

- **Solution**: Gradient clipping, advanced weight initialization, and skip connections help a computer learn things accurately and consistently.

- **Overfitting**
- Overfitting happens when a model knows too much about the training data, so it can't make good predictions about new data.
- As a result, the model performs well on the training data but struggles to make accurate predictions on new, unseen data.
- It's essential to address overfitting by employing techniques like regularization, cross-validation, and more diverse datasets to ensure the model generalizes well to unseen examples.
- **Solution**: Regularisation techniques help us ensure our models memorize the data and use what they've learned to make good predictions about new data. Techniques like dropout, L1/L2 regularisation, and early stopping can help us do this.

- **Data Augmentation and Preprocessing**

- Data augmentation and preprocessing are techniques used to provide better information to the model during training, enabling it to learn more effectively and make accurate predictions.

- **Solution**: Apply data augmentation techniques like rotation, translation, and flipping alongside data normalization and proper handling of missing values.

- **Label Noise**

- Training data sometimes need to be corrected, making it hard for computers to do things well.

- **Solution**: Using special kinds of math called "loss functions" can help ensure that the model you are using is not affected by label mistakes.

- **Imbalanced Datasets**

- Datasets can have too many of one type of thing and need more of another type. This can cause models not to work very well for things not represented as much.

- **Solution**: Classes can sometimes be uneven, meaning more people are in one group than another. To fix this, we can use special techniques like class weighting, oversampling, or data synthesis to ensure that all the classes have the same number of people.

- **Computational Resource Constraints**

- Training deep neural networks can be very difficult and take a lot of computer power, especially if the model is very big.

- **Solution**: Using multiple computers or special chips called GPUs and TPUs can help make learning faster and easier.

- **Hyperparameter Tuning**
- Deep neural networks have numerous hyperparameters that require careful tuning to achieve optimal performance.
- **Solution**: To efficiently find the best hyperparameters, such as Bayesian optimization or genetic algorithms, utilize automated hyperparameter optimization methods.
- **Convergence Speed**
- It is important to ensure a model works quickly when using lots of data and complicated designs.
- **Solution**: Adopt learning rate scheduling or adaptive algorithms like Adam or RMSprop to expedite convergence.

- **Activation Function Selection**

- Using the proper activation function when building a machine-learning model is important. This helps ensure the model works properly and yields correct results.

- **Solution**: ReLU and its variants (Leaky ReLU, Parametric ReLU) are popular choices due to their ability to mitigate vanishing activation issues.

- **Gradient Descent Optimization**

- Gradient descent algorithms help computers solve problems but sometimes need help when it is very difficult.

- **Solution**: Advanced techniques can help us navigate difficult problems better. Examples are stochastic gradient descent with momentum and Nesterov Accelerated Gradient.

- **Memory Constraints**
- Computers need a lot of memory to train large models and datasets, but they can work properly if there is enough memory.
- **Solution**: Reduce memory usage by applying model quantization, using mixed-precision training, or employing memory-efficient architectures like MobileNet or EfficientNet.
- **Transfer Learning and Domain Adaptation**
- Deep learning networks need lots of data to work well. If they don't get enough data or the data is different, they won't work as well.
- **Solution**: Leverage transfer learning or domain adaptation techniques to transfer knowledge from pre-trained models or related domains

- **Exploring Architecture Design Space**

- Designing buildings is difficult because there are many different ways to do it. Choosing the best way to create a building for a specific purpose can take time and effort.

- **Solution**: Use automated neural architecture search (NAS) algorithms to explore the design space and discover architectures tailored to the task.

- **Adversarial Attacks**

- Deep neural networks are unique ways of understanding data. But they can be tricked by minimal changes that we can't see. This can make them give wrong answers.

- **Solution**: Employ adversarial training, defensive distillation, or certified robustness methods to enhance the model's robustness against adversarial attacks.

- **Interpretability and Explainability**
- Understanding the decisions made by deep neural networks is crucial in critical applications like healthcare and autonomous driving.
- **Solution**: Adopt techniques such as LIME (Local Interpretable Model-Agnostic Explanations) or SHAP (SHapley Additive exPlanations) to explain model predictions.
- **Handling Sequential Data**
- Training deep neural networks on sequential data, such as time series or natural language sequences, presents unique challenges.
- **Solution**: Utilize specialized architectures like recurrent neural networks (RNNs) or transformers to handle sequential data effectively.

- **Limited Data**
- Training deep neural networks with limited labeled data is a common challenge, especially in specialized domains.
- **Solution**: Consider semi-supervised, transfer, or active learning to make the most of available data.
- **Catastrophic Forgetting**
- When a model forgets previously learned knowledge after training on new data, it encounters the issue of catastrophic forgetting.
- **Solution**: Implement techniques like elastic weight consolidation (EWC) or knowledge distillation to retain old knowledge during continual learning.

- **Hardware and Deployment Constraints**
- Using trained models on devices with not much computing power can be hard.
- **Solution**: Scientists use special techniques to make computer models run better on devices with limited resources.
- **Data Privacy and Security**
- When training computers to do complex tasks, it is essential to keep data private and ensure the computers are secure.
- **Solution**: Employ federated learning, secure aggregation, or differential privacy techniques to protect data and model privacy.

- **Long Training Times**

- Training deep neural networks is like doing a challenging puzzle. It takes a lot of time to assemble the puzzle, especially if it is vast and has a lot of pieces.

- **Solution**: Special tools like GPUs or TPUs can help us train our computers faster. We can also try using different computers simultaneously to make the training even quicker.

- **Exploding Memory Usage**

- Some models are too big and need a lot of space, so they are hard to use on regular computers.

- **Solution**: Explore memory-efficient architectures, use gradient checkpointing, or consider model parallelism for training.

- **Learning Rate Scheduling :**

- Setting an appropriate learning rate schedule can be challenging, affecting model convergence and performance.

- **Solution:** Using special learning rate schedules can help make learning easier and faster. These schedules can be used to help teach things in a better way.

- **Avoiding Local Minima**

- Deep neural networks can get stuck in local minima during training, impacting the model's final performance.

- **Solution**: Using unique strategies like simulated annealing, momentum-based optimization, and evolutionary algorithms can help us escape difficult spots.

- **Unstable Loss Surfaces**

- Finding the best way to do something can be very hard when there are many different options because the surface it is on is complicated and bumpy.

- **Solution**: Utilize weight noise injection, curvature-based optimization, or geometric methods to stabilize loss surfaces.

# Greedy Layer Wise Pre-Training

- Artificial intelligence has undergone a revolution thanks to neural networks, which have made significant strides possible in a number of areas like speech recognition, computer vision, and natural language processing. Deep neural network training, however, may be difficult, particularly when working with big, complicated datasets.

- One method that tackles some of these issues is greedy layer-wise pre-training, which initializes deep neural network settings layer by layer.

- Greedy layer-wise pre-training is used to initialize the parameters of deep neural networks layer by layer, beginning with the first layer and working through each one that follows.

- A layer is trained as if it were a stand-alone model at each step, using input from the layer before it and output to go to the layer after it. Typically, developing usable representations of the input data is the training aim.

# Processes of Greedy Layer-Wise Pre-Training

- The process of greedy layer-wise pre-training can be staged as follows:

- **Initialization:** The neural network's first layer is trained on its own using auto encoders and other unsupervised learning strategies. Learning a collection of features that highlight important elements of the input data is the aim.

- **Extracting Feature:** The activations of the first layer are utilized as features to train the subsequent layer after it has been trained. Each layer learns to represent the traits discovered by the layer before it in a higher-level abstraction when this process is repeated repeatedly.

- **Fine-Tuning:** The network is adjusted as a whole using supervised learning methods once every layer has been pre trained in this way. To maximize performance on a particular job, this entails simultaneously modifying all of the network's parameters using a labeled dataset.

# Advantages of Greedy Layer Wise Pre-Training

- Here are some of the advantages of Greedy Layer Wise Pre-Training:
- **Feature Learning and Representation:** At various degrees of abstraction, each layer of the network gains the ability to identify and extract pertinent characteristics from the incoming data. Pre-training is unsupervised, so the model may identify underlying structures and patterns in the data without needing labeled annotations. Consequently, the acquired representations often exhibit more information content and generalizability, resulting in enhanced performance on subsequent supervised tasks.
- **Regularization and Generalization:** Greedy layer-wise pre-training forces the model to acquire meaningful representations of the input data, which functions as a kind of regularization. By acting as a kind of regularization, the pre-trained weights direct the learning process to areas of the parameter space where there is a higher chance of good generalization to new data. This aids in avoiding overfitting, particularly in situations when training data is scarce.

- **Transfer Learning and Adaptability:** Greedy layer-wise pre-training makes it easier for a pre-trained model to transfer to new tasks or domains with little further training. This is known as transfer learning. The model is able to effectively adapt to new contexts and achieve acceptable performance even with insufficient labeled data, because of the learned features that capture general patterns in the data that are frequently transferable across other tasks or datasets.

- **Efficient Training Process:** Training every layer independently makes the entire process more effective and less prone to convergence problems. Later, the entire network may be fine-tuned using supervised learning. Pre-trained weights offer an excellent starting point for further training, which expedites the training process by lowering the number of iterations needed for convergence.

# Disadvantages of Greedy Layer Wise Pre-Training

- Greedy Layer Wise Pre-Training has various advantages but does come up with some limitations. Here are some of the disadvantages of Greedy Layer Wise Pre-Training:

- **Complexity and Training Time:** Greedy layer-wise pre-training teaches the neural network's layers independently using unsupervised learning, then uses supervised learning to fine-tune the entire network. This procedure can be costly and time-consuming in terms of processing, particularly for large-scale datasets and complex designs. Sequentially training more layers calls for more processing power and might not perform well for really deep networks.

- **Difficulty in Implementation:** It can be difficult to implement greedy layer-wise pre-training, especially in deep systems with several layers. Careful design and implementation are needed to ensure compatibility with the following fine-tuning processes, manage the transfer of pre-trained weights between layers, and coordinate the training process for each layer. The adoption of layer-wise pre-training may be hampered by its complexity, particularly for practitioners with little background in deep learning.

- **Dependency on Data Availability:** For unsupervised learning, greedy layer-wise pre-training needs access to a lot of unlabeled data. While this might not be a big deal for some domains or datasets, it might be problematic in situations when there is a lot of labeled data available but little or expensive unlabeled data to get. Other pre-training approaches or data augmentation methods could be more appropriate in certain circumstances.

- **Greedy Layer-Wise Pre-Training** is a training technique commonly used in the context of deep neural networks, particularly **unsupervised pre-training** of stacked layers, before fine-tuning the network on a supervised task.

- It is a **layer-by-layer approach** where each layer is trained independently in a greedy fashion, focusing on learning representations for the data at each layer of the network.

- This approach is often associated with **Deep Belief Networks (DBNs)** and **autoencoders**, and it was popularized as a way to address challenges in training deep neural networks, such as the **vanishing gradient problem** and poor initialization.

# Why Greedy Layer-Wise Pre-Training?

- Before **deep learning** achieved its current popularity with powerful optimization techniques, training deep networks from scratch was a very challenging task due to:

- **Vanishing Gradient Problem**: When training deep networks, gradients can become very small as they are propagated backward through many layers, making it hard for earlier layers to learn.

- **Bad Initialization**: Networks with many layers are sensitive to initial weights, and bad initialization can cause the model to get stuck in poor local minima.

- **Greedy layer-wise pre-training** helps to address these issues by training each layer individually and gradually learning good feature representations before combining them into a full network for the final task.

# Step-by-Step Process of Greedy Layer-Wise Pre-Training:

- **Unsupervised Training of the First Layer**:
  - The first layer is trained using unsupervised learning techniques, such as **autoencoders** or **Restricted Boltzmann Machines (RBMs)**. The goal here is to learn a good representation of the input data at the first layer.
  - The parameters (weights) of the first layer are initialized and updated using the data.

- **Freezing the First Layer**:
  - After the first layer is trained, its parameters are frozen (i.e., they are not updated during subsequent training of other layers).
  - The first layer's activations (or learned features) are now used as inputs for the second layer.

- **Training the Second Layer**:
  - The second layer is trained in the same manner as the first, but now it uses the activations of the first layer as its input.
  - Once the second layer is trained, its weights are frozen.

- **Repeat for All Layers**:
  - The same procedure is applied iteratively to all layers in the network. Each layer is trained independently, using the activations from the previous layer as its input.
  - After each layer is trained, the weights of the previous layers are frozen.

- **Fine-Tuning (Supervised Training)**:
  - Once all layers have been pre-trained, the full network is fine-tuned in a supervised manner using labeled data.
  - During fine-tuning, all the parameters of the network (from all layers) are updated to minimize the supervised loss, usually using **backpropagation**.

# Advantages of Greedy Layer-Wise Pre-Training:

- **Improved Initialization**:
  - Since each layer is trained independently and progressively, the network's weights are initialized in a way that captures good feature representations at each layer.
  - The approach avoids bad initialization, which often leads to slow convergence or poor local minima when training deep networks.
- **Avoiding the Vanishing Gradient Problem**:
  - By training each layer independently, the model learns to extract useful features from the data without facing the issues of very small gradients. This is particularly helpful for deep networks.
  - The individual layer-wise training can help with better weight updates in the earlier layers of the network.
- **Better Representations**:
  - Greedy layer-wise pre-training helps the model to progressively learn higher-level representations. The first layer might learn basic features like edges or textures, while deeper layers learn more abstract features (e.g., object parts, shapes).
- **Improved Performance**:
  - When fine-tuned with supervised data, the overall network tends to converge faster and often achieves better performance, especially in tasks where labeled data is limited.
- **Unsupervised Learning in Pre-Training**:
  - This method can take advantage of **unsupervised learning** techniques (like autoencoders or RBMs) to learn representations from unlabeled data before using labeled data in the final step.

# Intuitive perspective of optimization

- The goal of the optimization of *DNN* is to find the best parameters w to minimize the loss function *f(w, x, y)*, using *f(w)* bellow for simplification, subject to *x, y*, where *x* are the data and *y* are the labels.

- The gradient descent (*GD*) is the <span style="color:red">most frequently used optimization method for machine learning</span>. In this method, we need another parameter called learning rate, $\alpha$.

Now we start the process of gradient descent, at each batch/step *t:*

1. Calculate the gradient of each parameter: $g_t = \nabla f(w_t)$

2. Calculate the first order and/or second order momentum of gradient:
   $m_t = \phi(g_1, g_2, ..., g_t)$ and $V_t = \psi(g_1^2, g_2^2, ..., g_t^2)$

3. Calculate the update value for w: $\eta_t = \alpha * m_t / (\sqrt{V_T} + \epsilon)$

4. Update w: $w_{t+1} = w_t - \eta_t$

- Gradient descent is a process of iterative optimization, where you gradually move toward the minimum of a function (like a valley) by following the steepest downward slope, adjusting your steps based on the gradient (slope) at each point.
- By fine-tuning the step size (learning rate), gradient descent helps you find the best possible solution to a problem, often used in training machine learning models.

# Why Gradient Descent?

- **Gradient Descent** (GD) is the most widely used optimization algorithm for training machine learning models, particularly in problems involving continuous optimization, such as training neural networks, linear regression, and logistic regression.

- The basic idea behind gradient descent is simple: <span style="color:red">**minimize the loss (or cost) function** of the model by iteratively updating the model's parameters in the direction that reduces the loss.</span>

- In machine learning, <span style="color:red">the goal is often to minimize the **loss function** (or cost function), which measures how far the model's predictions are from the actual values</span>. <span style="color:purple">**Gradient descent helps find the parameter values that lead to the minimum loss**</span>.

# *Batch gradient descent (BGD) \*\*\**

- *BGD* calculates the gradients with entire training dataset for only one update of parameters. It can be very slow to converge.

- If the training dataset is too large to fill into the memory, *BGD* becomes intractable. Moreover, *BGD* is not compatible to update a model online, for example, with new data on-the-fly.

- Batch Gradient Descent is a reliable and straightforward method for optimizing a function. It's great for smaller datasets or situations where computational cost is not a concern.

- However, for large-scale problems, alternatives like **Stochastic Gradient Descent (SGD)** or **Mini-Batch Gradient Descent** are often preferred due to their faster convergence and lower computational cost.

# *Stochastic gradient descent (*SGD*) \*\*\**

- *SGD* calculates the gradients with one data. The calculation becomes faster, but the process of gradient descent becomes fluctuating. The direction of gradient descent calculated from only one data is not globally stable. It can be even opposite to the real gradient descent direction.

# Key Characteristics of Stochastic Gradient Descent (SGD):****

- **Faster Updates**:
  - **SGD** updates the model parameters after processing each individual data point, so it can converge more quickly, especially for large datasets.

- **Noisy Convergence**:
  - Because only one data point is used to compute the gradient, the updates are more erratic, causing the algorithm to oscillate around the minimum, rather than smoothly converging.
  - This noise can sometimes help escape local minima (in non-convex problems like neural networks), but it can also cause the algorithm to take longer to settle at the global minimum.

- **More Frequent Updates**:
  - Since it processes one data point at a time, **SGD** performs updates much more frequently than Batch Gradient Descent, which requires processing the entire dataset for each update.

- **Memory Efficient**:
  - Since SGD only needs a single data point to compute the gradient, it doesn't require the entire dataset to be stored in memory, which makes it much more memory-efficient for large datasets.

- **Learning Rate Sensitivity**:
  - The noisier updates can make SGD sensitive to the learning rate $\alpha$\alpha$\alpha$. If the learning rate is too large, the parameter updates might overshoot the minimum. If it's too small, the algorithm may take a long time to converge.

- Stochastic Gradient Descent (SGD) is an optimization algorithm that works by updating the model parameters after processing each individual data point.

- It is faster and more memory-efficient than **Batch Gradient Descent**, especially for large datasets.

- However, the updates are noisier and more erratic, which can cause slower convergence or oscillations around the minimum.

- Variants of **SGD**, like **Mini-Batch Gradient Descent** or methods with momentum and adaptive learning rates, can help mitigate these issues and improve convergence stability.

# Adagrad

- **Adagrad** (short for Adaptive Gradient Algorithm) is an optimization algorithm designed to adapt the learning rate for each parameter during the optimization process.

- Unlike traditional gradient descent methods, which use a fixed learning rate for all parameters, **Adagrad** adjusts the learning rate based on the historical gradients of each parameter.

- The goal is to ensure that the model learns efficiently, especially for sparse data or highly variable gradients across parameters.

# Key Concept of Adagrad:

- **Adaptive Learning Rate**: **Adagrad** adapts the learning rate for each parameter based on the past gradients. Parameters that have large gradients in the past receive smaller learning rates, and parameters with smaller gradients receive larger learning rates.

- **Scaling of Updates**: The algorithm scales the updates to each parameter inversely proportional to the square root of the accumulated gradient information over time. This helps prevent large updates for frequently updated parameters and smaller updates for infrequently updated ones.

# Advantages of Adagrad:

- **Adapts to Sparse Data**:
  - **Adagrad** is particularly useful for problems where the dataset has sparse features (e.g., text data, with many zero entries). It automatically adjusts learning rates to perform better on sparse data.

- **No Need for Learning Rate Tuning**:
  - One of the main advantages of **Adagrad** is that it does not require manual tuning of the learning rate, as the algorithm adapts the learning rate for each parameter.

- **Improves Efficiency**:
  - Since it adapts the learning rates based on past gradients, **Adagrad** often leads to faster convergence, especially for parameters that need larger updates initially.

- **Works Well with Sparse Gradients**:
  - **Adagrad** is particularly effective when working with sparse gradients, as it helps in boosting the learning rate for infrequent updates and decelerates for frequent ones.

# RMSprop

- **RMSprop** (Root Mean Square Propagation) is an adaptive learning rate optimization algorithm that addresses some of the issues with **Adagrad**. Specifically, **RMSprop** modifies **Adagrad** by reducing the problem of a learning rate that decays too quickly, allowing for more effective training over a longer period of time.

- Like **Adagrad**, **RMSprop** adapts the learning rate for each parameter. However, **RMSprop** improves on **Adagrad** by using an exponentially decaying average of past gradients, rather than summing the squared gradients over all time. This ensures that the learning rate doesn't shrink too rapidly, allowing the algorithm to continue learning effectively throughout the training process.

# Intuition Behind RMSprop:

- **Decay of Past Gradients**: Unlike **Adagrad**, where the learning rate is updated based on the cumulative sum of squared gradients, **RMSprop** uses an exponentially weighted average of squared gradients. This prevents the learning rate from decaying too rapidly and helps maintain reasonable step sizes as training continues.

- **Adaptive Learning Rate**: Just like **Adagrad**, **RMSprop** adapts the learning rate for each parameter, but it balances the influence of past gradients. This makes it more stable and effective than **Adagrad** for many practical problems.

# Advantages of RMSprop:

- **Prevents Rapid Learning Rate Decay**:
  - **RMSprop** effectively <span style="color:red">mitigates the problem of a learning rate that decays too quickly</span>, as seen in **Adagrad**, by using an <span style="color:red">exponentially decaying average</span> of the squared gradients.

- **Faster Convergence**:
  - The adaptive learning rate helps speed up convergence, especially in cases with noisy or sparse gradients, such as in deep learning problems.

- **Effective for Non-Stationary Objectives**:
  - **RMSprop** is well-suited for problems where the data or the objective function is non-stationary (e.g., online learning or time series), because it adapts the learning rate based on the current gradients.

- **Memory Efficient**:
  - Since it <span style="color:red">only requires storing the exponentially decaying averages of squared gradients</span> (<span style="color:red">instead of maintaining all past gradients as in **Adagrad**</span>), it is more memory efficient.

- **Works Well for Neural Networks**:
  - **RMSprop** is often used in the training of deep neural networks, as it performs well in the presence of <span style="color:red">sparse gradients or gradients that have different magnitudes across</span> parameters.

# Adam (Adaptive Moment Estimation)

- **Adam** is a popular optimization algorithm used to train machine learning models, particularly in deep learning.

- It combines the benefits of two other extensions of gradient descent: **Momentum** and **RMSprop** (Root Mean Square Propagation).

- Adam has become widely adopted because of its adaptive learning rates, efficiency, and ability to handle sparse gradients and noisy data.

# Why Adam?

- Training deep neural networks can be challenging, especially when dealing with large datasets and complex models.

-  The **gradient descent** algorithm, while effective, often requires manual tuning of hyperparameters like the **learning rate**.

- If the learning rate is not chosen carefully, the training can either be too slow or fail to converge properly.

- Adam automates the process of adjusting the learning rate during training, which results in faster convergence, better stability, and less manual tuning.

# How Adam Works

- Adam is an adaptive learning rate optimization algorithm that computes individual learning rates for each parameter based on estimates of both **first moment** (mean) and **second moment** (uncentered variance) of the gradients. These moment estimates are used to adjust the learning rate of each parameter, making Adam both efficient and effective for a wide range of problems.

- The key components of Adam are:

- **Momentum (First Moment Estimate)**: It helps the optimization algorithm maintain a smooth trajectory towards the optimum by incorporating the past gradients. This is similar to the **Momentum** optimization method.

- **RMSprop (Second Moment Estimate)**: RMSprop uses a moving average of the squared gradients to adapt the learning rate. This helps the optimization process be more stable by adjusting the learning rate based on the magnitude of the gradients.

# Advantages of Adam:

- **Adaptive Learning Rates**: Adam adapts the learning rate for each parameter, allowing for efficient training, particularly when different parameters require different rates of updates.

- **Combines Momentum and RMSprop**: Adam combines the benefits of both **Momentum** (which helps smooth out the gradients and accelerates convergence) and **RMSprop** (which helps adjust the learning rate for each parameter individually).

- **Efficient Computation**: The algorithm is computationally efficient and requires little memory overhead, making it suitable for large datasets and deep networks.

- **Bias Correction**: The bias-correction terms ensure that the moment estimates are accurate even in the early stages of training when they might be biased towards zero.

- **Good Default Choice**: Adam has been found to work well for a wide range of machine learning tasks out-of-the-box, making it a popular choice for deep learning applications.

- **Adam** has become one of the most popular and effective optimization algorithms for deep learning due to its adaptive learning rates, efficient computation, and ability to handle sparse gradients and noisy data.

- It combines the advantages of both **Momentum** and **RMSprop**, making it suitable for a wide variety of machine learning tasks. While it may require some hyperparameter tuning, Adam generally works well with its default settings and has proven to be robust in a wide range of applications.