# A Comprehensive Guide to Optimizing ConvNeXt Model Training

## Foundational Hyperparameters and Regularization Strategies

The selection of foundational hyperparameters and regularization techniques forms the bedrock of any successful deep learning experiment, particularly for complex architectures like ConvNeXt. The provided research landscape reveals a strong consensus on several core components, while also highlighting areas where strategic variation can yield significant performance gains. This section synthesizes these findings into a coherent framework for establishing a robust starting point for your hyperparameter sweep. The primary focus will be on optimizer choice, learning rate schedules, and the application of various regularization mechanisms such as stochastic depth, label smoothing, and gradient clipping. These elements are not merely technical details but are deeply intertwined with the optimization dynamics of modern convolutional networks, influencing everything from initial convergence speed to final generalization performance.

A dominant theme emerging from recent literature is the near-universal adoption of the AdamW optimizer for training ConvNeXt models [3] [4] [14] [15] [18] [20] [21] [23]. While the standard Adam optimizer is also employed successfully in numerous studies [3] [15] [18] [21], AdamW has become the de facto standard for transformer-style architectures, which ConvNeXt emulates with its block structure [10]. The critical distinction lies in how weight decay is applied; AdamW incorporates it explicitly within the optimization step, whereas older implementations often applied it separately, leading to suboptimal results. This explicit regularization helps prevent overfitting and promotes better generalization, making it a superior choice for state-of-the-art models. For both Adam and AdamW, the momentum parameters, β1 and β2, are consistently set to the widely accepted defaults of 0.9 and 0.999, respectively [15] [18] [23]. The initial learning rate (LR) for Adam-family optimizers exhibits a clustering around two primary values: $1 \times 10^{-4}$ and $1 \times 10^{-3}$. Several sources report using an initial LR of $1 \times 10^{-4}$ for tasks ranging from nuclear segmentation to image classification [3] [14] [15] [20]. Others have found success with a higher initial LR of

$1\times10^{-3}$, particularly in generative and captioning tasks [3] [23] [33]. This suggests that a broad initial sweep range of at least $1\times10^{-5}$ to $1\times10^{-2}$ is warranted to explore this landscape effectively. The weight decay parameter for AdamW shows similar consistency, with values typically falling between $1\times10^{-4}$ and $1\times10^{-2}$. Commonly cited values include 0.01 [3] [14] [15] [20] and 0.0001 [3] [18]. A systematic sweep exploring values from $1\times10^{-5}$ to $1\times10^{-1}$ would provide a thorough exploration of this crucial regularization hyperparameter.

Beyond the optimizer itself, the learning rate schedule plays a pivotal role in navigating the complex loss landscape of deep neural networks. The evidence strongly indicates that adaptive schedules are far more effective than a simple constant learning rate throughout training. Cosine annealing is a frequently used and highly effective strategy, often implemented with a warmup period to stabilize the early stages of training before gradually reducing the LR towards the end [3] [14] [15] [20] [21]. One study on malaria blood smear classification specifically utilized a `OneCycleLR` scheduler, which is a form of cyclical scheduling that can offer benefits in certain scenarios [4]. The duration of the warmup period varies, with some studies employing a short warmup of just 5 epochs [14] [20] and others using a more extended warmup of up to 20 epochs [23]. This variability suggests that the optimal warmup length may depend on factors such as dataset size, model complexity, and the chosen optimizer. Therefore, including a warmup phase in any hyperparameter sweep is a sound practice, with the duration itself becoming a tunable parameter. The combination of a well-chosen optimizer, a carefully scheduled learning rate, and appropriate weight decay creates a powerful foundation for stable and efficient training.

Regularization extends beyond the optimizer and encompasses a variety of techniques designed to improve a model's ability to generalize to unseen data. For ConvNeXt models, one of the most important architectural components is stochastic depth, also known as DropPath [10]. This technique randomly drops entire blocks during training, effectively creating an ensemble of shallower networks. It is not merely an optional hyperparameter to be tuned in every experiment but a core feature of the ConvNeXt architecture. The rate at which blocks are dropped, `drop_path_rate`, is configurable. In many implementations, it is set to a small value, such as 0.1 [11] [23], while in others, it is left at the default of 0.0 [10]. For very large models, experimenting with different `drop_path_rate` values could be a fruitful avenue of research to understand its impact on convergence and performance. Another effective regularization technique demonstrated in the literature is label smoothing. In a successful malaria detection task, applying label

smoothing with a smoothing factor α of 0.1 resulted in a 98.1% accuracy score, suggesting it can help mitigate overconfidence in predictions and improve robustness [4] . Finally, gradient clipping is a crucial technique for preventing training instability, especially in generative and segmentation tasks. Thresholds of 0.1 [3] and 1.0 [15] have been shown to be effective in controlling exploding gradients without hindering convergence.

The table below summarizes the foundational hyperparameters and regularizers discussed, providing a synthesis of values reported across multiple research contexts.

| Hyperparameter/ Technique | Recommended Value / Range | Rationale and Supporting Evidence |
|---|---|---|
| Optimizer | AdamW | Considered the modern standard for transformer-style architectures, providing explicit weight decay regularization. [3] [4] [14] [15] [18] [20] [21] [23] |
| | Adam | Also proven effective, but AdamW is generally preferred for state-of-the-art results. [3] [15] [18] [21] |
| Initial Learning Rate (LR) | $1 \times 10^{-4}$ to $1 \times 10^{-3}$ | Initial LRs cluster around these values. A broader sweep from $1 \times 10^{-5}$ to $1 \times 10^{-2}$ is recommended. [3] [14] [15] [20] [23] [33] |
| Weight Decay | $1 \times 10^{-4}$ to $1 \times 10^{-2}$ | Values consistently fall in this range. Sweep from $1 \times 10^{-5}$ to $1 \times 10^{-1}$ for thorough exploration. [3] [14] [15] [18] [20] |
| β1 (for Adam/W) | 0.9 | Standard default for Adam-family optimizers. [15] [18] [23] |
| β2 (for Adam/W) | 0.999 | Standard default for Adam-family optimizers. [15] [18] [23] |
| Learning Rate Schedule | Cosine Annealing | Highly effective, often combined with a warmup phase to stabilize training. [3] [14] [15] [20] [21] |
| | OneCycleLR | Proven effective for achieving high accuracy in a limited number of epochs. [4] |
| Warmup Period | 5 - 20 epochs | Duration varies by task and model; serves to stabilize early training. [14] [20] [23] |
| DropPath Rate (`drop_path_rate`) | 0.0 - 0.3 | Core architectural component of ConvNeXt. Default is often 0.0, but values up to 0.3 have been explored. [10] [11] [23] |
| Label Smoothing | α = 0.1 | Can improve generalization and model confidence calibration. [4] |
| Gradient Clipping | Threshold of 0.1 - 1.0 | Prevents exploding gradients and improves training stability. [3] [15] |

In conclusion, establishing a solid baseline requires careful consideration of these interconnected hyperparameters. The consensus leans heavily towards using AdamW with a moderate initial learning rate ($1 \times 10^{-4}$), a small weight decay ($1 \times 10^{-2}$), and a cosine annealing schedule with a warmup period. This combination provides a robust and stable training process that serves as an

excellent foundation upon which more advanced techniques and architectural variations can be built. The subsequent sections will build upon this foundation, addressing the practicalities of scaling these configurations to large-scale hardware and leveraging modern tooling to maximize efficiency.

# Batch Size Scaling and Mixed-Precision Training Best Practices

Effective utilization of modern GPU hardware, such as the NVIDIA RTX 4090, hinges critically on two interrelated concepts: intelligent batch size scaling and the widespread adoption of mixed-precision training. The RTX 4090's Tensor Cores are engineered to deliver substantial performance boosts for matrix operations performed in lower precision formats like FP16, making mixed-precision not just an option but a necessity for competitive training times [13] [22]. Simultaneously, managing batch size—particularly across multiple GPUs—is essential for fully saturating the computational capacity of the hardware and achieving high throughput. This section delves into the principles of scaling batch size proportionally with learning rates, the mechanics and implementation of automatic mixed-precision (AMP), and the profound impact these choices have on both training speed and numerical stability.

The principle of scaling the learning rate proportionally with the batch size is a cornerstone of deep learning best practices, and it is explicitly mentioned as a guideline for adapting ConvNeXt models [9]. This heuristic is rooted in the idea that a larger batch size provides a more accurate estimate of the true gradient, thus allowing for a larger step size (i.e., a higher learning rate) without causing divergence. However, a more nuanced understanding of adaptive optimizers like AdamW reveals a non-monotonic relationship between batch size and optimal learning rate [28]. Instead of a simple linear increase, there exists a "surge peak" batch size, denoted B*, where the optimal learning rate reaches its maximum. Beyond this point, further increases in batch size cause the optimal learning rate to decrease. This implies that simply scaling the LR linearly with the batch size may not always lead to the fastest convergence. For your research, this suggests that a hyperparameter sweep should not only test a range of batch sizes but also systematically explore the corresponding optimal learning rate for each, paying close attention to this potential "surge" phenomenon. This could uncover batch-

size/LR pairs that achieve faster training than those found via simple proportional scaling alone.

Mixed-precision training is the standard for efficient training on modern NVIDIA GPUs, and its implementation has become remarkably straightforward with PyTorch 2.x [1] . The core concept involves performing most of the forward and backward passes in the 16-bit floating-point format (FP16 or BF16) while maintaining a master copy of the model's weights in the 32-bit format (FP32) [13] . This hybrid approach offers two major advantages. First, it significantly reduces the memory footprint, as storing activations and gradients in FP16 instead of FP32 cuts memory usage by half [22] . Second, and more importantly, it dramatically accelerates computation. The RTX 4090's Tensor Cores are optimized to execute FP16 matrix multiplication and accumulation operations much faster than their FP32 counterparts, leading to substantial throughput improvements [22] . Benchmarks on other NVIDIA hardware show consistent speedups; for instance, an A100 GPU achieved a 1.72x speedup for ResNet50 in FP16 compared to FP32 [24] . To maintain numerical stability in this low-precision environment, it is critical to use a technique called loss scaling. Gradients computed in FP16 can become so small that they underflow to zero. Loss scaling involves multiplying the final loss by a large scalar value before backpropagation, ensuring that the resulting gradients are large enough to be represented accurately in FP16. After computing the gradients, they are then unscaled before updating the FP32 master weights. PyTorch's Automatic Mixed Precision (AMP) library, specifically `torch.cuda.amp.autocast` and `torch.cuda.amp.GradScaler`, automates this entire process, making it accessible out-of-the-box in modern environments [1] [4] [13] . This seamless integration removes a significant barrier to adopting mixed-precision, which is now considered a fundamental requirement for efficient ConvNeXt training on consumer-grade hardware like the RTX 4090.

The practical application of these principles in a multi-GPU setting requires careful management of batch distribution. The total effective batch size is the product of the per-GPU batch size, the number of GPUs, and the number of gradient accumulation steps. A clear example of this is seen in the training of InceptionNeXt-T, where a total batch size of 4096 was achieved on 8 GPUs by using a per-GPU batch size of 128 and 4 gradient accumulation steps $(4096=8 \text{ GPUs} \times 128 \text{ per-GPU batch size} \times 4 \text{ accumulation steps})$ [23] . Similarly, a ConvNeXt-based nuclear segmentation model was trained on 4 RTX 4090s with a per-GPU batch size of 16, resulting in a total batch size of 64 $(4 \text{ GPUs} \times 16 \text{ per-GPU batch size} = 64)$ [14] [20] . The choice of per-GPU batch size is

constrained by the available VRAM on each GPU. With the RTX 4090's 24GB of memory, you can typically accommodate larger per-GPU batch sizes than on cards with less memory. Your hyperparameter sweep should therefore explore a wide range of per-GPU batch sizes, from smaller values suitable for very deep models to larger values that can fully leverage the 24GB memory budget. For each per-GPU batch size, you should pair it with a corresponding scaled learning rate, guided by the principles discussed earlier.

The table below outlines a practical framework for configuring batch size and mixed-precision training for your ConvNeXt experiments.

| Configuration Aspect | Recommendation / Strategy | Implementation Details and Rationale |
| --- | --- | --- |
| Batch Size Scaling | Scale LR proportionally with batch size. | Follows established best practices to maintain stable training dynamics. [9] |
| | Explore the "surge peak" phenomenon. | Adaptive optimizers like AdamW have an optimal LR that peaks at a certain batch size before decreasing. [28] |
| Per-GPU Batch Size | Start with a moderate size (e.g., 16-32) and scale up. | Constrained by 24GB VRAM on RTX 4090. Larger batches generally lead to higher throughput. [14] [20] [23] |
| Total Batch Size | Aim for a large total batch size (e.g., 64-1024). | Maximizes GPU utilization and throughput. Total batch size = Per-GPU batch size × #GPUs × Accumulation Steps. [14] [20] [23] |
| Gradient Accumulation | Use when per-GPU batch size is limited by VRAM. | Allows for an effective large batch size without requiring prohibitively large VRAM per GPU. [2] [23] |
| Mixed-Precision (AMP) | Enable by default. | Reduces VRAM usage by ~50% and accelerates computation via Tensor Cores on RTX 4090. [1] [13] [22] |
| Implementation | Use `torch.cuda.amp.autocast` and `GradScaler`. | PyTorch 2.x provides a seamless, automated implementation that handles loss scaling. [1] [4] |
| Numerical Stability | Ensure `GradScaler` is used correctly. | Prevents gradient underflow/overflow in the FP16 environment, which is critical for stable training. [4] [13] |

In summary, optimizing batch size and embracing mixed-precision are not optional extras but fundamental requirements for efficient training of ConvNeXt models on the RTX 4090 platform. By systematically exploring the relationship between batch size, learning rate, and gradient accumulation, and by fully leveraging the computational power of Tensor Cores through AMP, you can construct a training pipeline that maximizes both speed and resource utilization. This approach provides a solid foundation for your research, freeing up computational resources to be focused on more complex aspects of the hyperparameter space.

# Multi-GPU Utilization and Distributed Training Paradigms

Scaling ConvNeXt training from a single GPU to a multi-GPU cluster, such as an 8x RTX 4090 setup, introduces a new layer of complexity centered on communication overhead, synchronization, and load balancing. The choice of distributed training paradigm is arguably the most critical decision for achieving high GPU utilization and minimizing training time. The provided research highlights a stark contrast in performance between different strategies, with modern approaches like DistributedDataParallel (DDP) offering significant advantages over older, simpler methods. This section provides a detailed analysis of these paradigms, practical guidelines for configuring them on your hardware, and system-level tuning tips to ensure stable and efficient multi-GPU operation.

The empirical evidence overwhelmingly points to DistributedDataParallel (DDP) as the superior choice for multi-GPU training of ConvNeXt models. A benchmark conducted on an 8-GPU node revealed that DDP achieved approximately 57% average GPU utilization and completed training in about 190 seconds [25] . This performance stood in stark contrast to other advanced methods like FairScale FSDP (~60% utilization, 203s), native FSDP (~53% utilization, 210s), and DeepSpeed (~45% utilization, 1498s), and was vastly superior to the deprecated DataParallel, which only managed ~13% utilization and took 648 seconds [25] . The reason for DDP's effectiveness is its architecture: it creates a separate process for each GPU, with each process holding a complete replica of the model. During the forward pass, data is partitioned and fed to each process. The gradients are then synchronized across all processes after the backward pass, ensuring that all model replicas remain identical. This "all-reduce" synchronization step is handled efficiently by the NCCL backend, which is optimized for high-speed inter-GPU communication on NVIDIA hardware [25] . For your research, starting with DDP is a highly recommended and likely optimal strategy for all multi-GPU experiments. The provided Runpod environment, which includes the RTX 4090, is explicitly validated for computer vision tasks and supports multi-GPU pod deployments, making it a suitable platform for implementing DDP-based training [1] .

While DataParallel is no longer recommended due to its poor performance and scalability limitations [25] , it is worth mentioning for historical context. DataParallel operates by wrapping the model in a single process and replicating the data across GPUs. Each GPU receives a slice of the data, performs a forward and backward pass independently, and then all gradients are collected and averaged in the main

process before the optimizer update. This central bottleneck makes it inefficient, especially as the number of GPUs increases, explaining its abysmal performance in benchmarks [25] . Fully Sharded Data Parallel (FSDP) and DeepSpeed are more advanced strategies designed to handle models that are too large to fit entirely in the memory of a single GPU. They work by sharding the model's parameters, gradients, and optimizer states across all available GPUs. While powerful for ultra-large models, they introduce additional communication overhead and complexity. As noted in the benchmark, their performance on smaller models can be inferior to DDP [25] . Given that your model is "very large," it is possible that FSDP might eventually become necessary. However, for initial sweeps and even potentially for your largest models, DDP remains a strong contender. If you encounter OOM errors with DDP, benchmarking FairScale FSDP against native PyTorch FSDP would be a logical next step, keeping in mind the performance trade-offs observed in the provided data.

Configuring a multi-GPU run involves determining the per-GPU batch size, the total batch size, and the number of gradient accumulation steps. As previously discussed, the total batch size is calculated as

$Total_BatchSize = PerGPU_BatchSize \times Number_of GPUs \times Gradient Accumulation Steps$

[23] . The goal is to find the largest possible total batch size that fits comfortably within the collective VRAM of all GPUs. For an 8x RTX 4090 setup, the collective VRAM is 192 GB, which allows for extremely large batch sizes. A good starting point for a sweep could involve varying the number of GPUs from 1 to 8 while keeping the per-GPU batch size constant, thereby scaling the total batch size linearly. Alternatively, you could fix the total batch size and vary the number of GPUs and accumulation steps accordingly. For example, to achieve a total batch size of 512, you could use:1 GPU with a batch size of 512.2 GPUs with a batch size of 256 each and 1 accumulation step.4 GPUs with a batch size of 128 each and 1 accumulation step.8 GPUs with a batch size of 64 each and 1 accumulation step.

This systematic approach allows you to measure both the raw speedup and the scaling efficiency (how close the actual speedup is to the ideal linear speedup) of your training pipeline.

Finally, ensuring stable operation in a multi-GPU environment often requires some system-level tuning. One common issue in multi-GPU setups, particularly when using AMP, is intermittent hangs during initialization or training. One source resolved such hangs on 4 GPUs by disabling cuDNN's auto-tuner (`torch.backends.cudnn.benchmark = False`) and enabling explicit CUDA synchronization (`torch.cuda.synchronize()`) before the main training loop [2] .

While this fix did not completely resolve all hangs, it highlights the importance of disabling features that can interfere with deterministic execution in a parallel environment. Another useful tip is to manage the number of CPU threads allocated to PyTorch via the `OMP_NUM_THREADS` environment variable. Setting this variable before importing PyTorch can help balance the workload between the CPUs and the GPUs, reducing stalls and improving overall throughput, especially in data loading pipelines [2] . These low-cost, high-impact adjustments should be part of your standard multi-GPU training script configuration.

The following table compares the primary distributed training paradigms discussed, summarizing their suitability for your ConvNeXt research.

| Paradigm | Average GPU Utilization (Benchmark) | Relative Speed (Benchmark) | Scalability | Key Advantage | Key Disadvantage |
|---|---|---|---|---|---|
| DataParallel | ~13% | Slowest (Baseline) | Poor | Simple to implement. | Centralized gradient aggregation bottleneck; very slow on >1 GPU. [25] |
| DistributedDataParallel (DDP) | ~57% | Fastest (~3x FSDP variants) | Excellent | High efficiency, low communication overhead, excellent scalability. | Slightly more complex setup than DataParallel. [25] |
| Native FSDP | ~53% | Slower than DDP (~1.1x slower) | Very Good | Scales to models larger than single-GPU VRAM. | Higher communication overhead than DDP for smaller models. [25] |
| FairScale FSDP | ~60% | Faster than DDP (~0.9x DDP) | Very Good | Similar to Native FSDP but with a different API. | Slight performance overhead compared to DDP. [25] |
| DeepSpeed | ~45% | Slowest (~3.2x DDP) | Excellent | Supports ZeRO optimization levels for extreme memory savings. | Highest communication overhead and complexity; not always fastest. [25] |

To conclude, achieving high GPU utilization for your ConvNeXt training on 1x and 8x RTX 4090s is primarily a function of choosing the right distributed training strategy. DDP emerges as the clear winner for efficiency and scalability. By combining DDP with a thoughtful strategy for batch size scaling and gradient accumulation, and by incorporating minor system-level optimizations, you can build a robust and highly performant multi-GPU training pipeline that fully leverages the power of your hardware.

# Advanced System-Level Optimizations and Tooling

Beyond the foundational hyperparameters and distributed training frameworks, a suite of advanced system-level optimizations and modern software tooling can significantly enhance the efficiency, stability, and speed of ConvNeXt training. These techniques operate at a lower level than traditional hyperparameter tuning, focusing on compiling code into highly efficient kernels, trading computation for memory, and leveraging specialized hardware features. Two of the most impactful technologies for this purpose are PyTorch's `torch.compile` and gradient checkpointing. Furthermore, understanding the nuances of normalization layers and their interaction with data preprocessing is crucial for avoiding subtle pitfalls that can degrade performance. This section explores these advanced topics, providing actionable guidelines for integrating them into your research workflow to push the boundaries of training performance on the RTX 4090 platform.

`torch.compile` is a revolutionary feature introduced in PyTorch 2.0 that JIT-compiles Python code into optimized, fused kernels, drastically reducing overhead and accelerating execution [5] . Its benefits are twofold: it minimizes Python interpreter overhead and reduces expensive GPU read/write operations by fusing operators together [5] . Empirical results demonstrate its power, showing a median speedup of 2.33x on a CUDA-optimized tensor operation [5] . This acceleration is orthogonal to hyperparameter tuning; it applies directly to the model's computation graph and can yield significant performance gains across the board. Crucially, `torch.compile` works natively with mixed-precision, seamlessly integrating with `torch.cuda.amp[[5]]`. This makes it an exceptionally powerful tool for ConvNeXt training on the RTX 4090, where kernel efficiency is paramount. An empirical evaluation confirmed that `torch.compile` provides significant performance benefits on ConvNeXt models running on an RTX 4090, closely approximating baseline accuracy while leveraging optimized kernels for inference [32] . However, there is a caveat: `torch.compile` can struggle with data-dependent control flow, such as `if` statements that depend on tensor values. When such constructs are present, the compiler may break the graph, leading to partial compilation and reduced optimization opportunities [5] . To achieve full-graph compilation, these conditional branches must be explicitly converted to functions compatible with `torch.cond[[5]]`. For a standard ConvNeXt architecture, which consists primarily of sequential layers and convolutions, full-graph compilation is likely achievable with minimal code changes. Therefore, a key recommendation for your research is to test `model.compile(fullgraph=True)` on your target model. Benchmarking

its impact on training time versus the overhead of the initial compilation phase is a valuable experiment in itself.

Gradient checkpointing is another powerful technique for managing memory consumption, particularly for very deep feed-forward networks like ConvNeXt [26]. The core idea is to selectively discard intermediate activations during the forward pass and re-compute them on-demand during the backward pass. This trades increased computation (an extra forward pass) for a significant reduction in memory usage. The theoretical memory saving can be dramatic, reducing the memory footprint from $O(n)$ to $O(sqrt(n))$ for an n-layer model [26]. This technique is invaluable when VRAM becomes the limiting factor, allowing you to train larger models or use larger batch sizes than would otherwise be possible. However, its practical application requires careful consideration. Research on its implementation in PyTorch shows that the marginal benefit in peak memory reduction diminishes quickly with the number of checkpoints. Using more segments (checkpoints) yields progressively smaller reductions in peak memory but incurs a greater runtime penalty due to the increasing number of recomputations [27]. For a ConvNeXt model, a sensible strategy would be to apply gradient checkpointing at the block level, i.e., wrap groups of ConvNeXt blocks in the `torch.utils.checkpoint` function. You should conduct a benchmark to determine the optimal number of segments (or blocks per segment) that provides the desired memory saving without introducing an unacceptable slowdown in training time. This makes gradient checkpointing a highly tunable hyperparameter in its own right.

The choice of normalization layers is a defining characteristic of the ConvNeXt architecture and has profound implications for training stability and performance. Unlike many previous CNNs that relied on Batch Normalization (BN), ConvNeXt employs Layer Normalization (LN) [17]. This shift was motivated by BN's sensitivity to batch size and its tendency to cause instability and overfitting in some scenarios. LN operates on the channel dimension within a single sample, making it independent of the batch size and often more stable. In an ablation study on cocoa disease detection, a ConvNeXt Tiny model with LN achieved an 88.14% F1 score, which was 2.84% higher than a ResNet18 with BN [17]. This demonstrates the tangible benefits of LN for this architecture. Interestingly, the interaction between normalization layers and data preprocessing can be complex. In the same cocoa disease detection study, removing the standard image normalization step (which subtracts mean and divides by standard deviation) actually improved all metrics and reduced overfitting [17]. This occurred because the image normalization process distorted the visibility of lesions in the cocoa pod images, counteracting the benefits of the LN layer. This

finding underscores a critical lesson: the choice of normalization and preprocessing techniques should be tailored to the specific characteristics of the dataset and the model architecture. For your research, it is advisable to experiment with different combinations of input normalization and the inherent normalization within the ConvNeXt blocks to see if a custom preprocessing pipeline can yield better results.

The table below summarizes these advanced optimization techniques and provides guidance on their application.

| Technique | Primary Benefit | Implementation Guidance | Potential Trade-off |
|---|---|---|---|
| `torch.compile` | Significant speedup (up to 2.33x) by fusing kernels and reducing overhead. | Use `model.compile(fullgraph=True)` for best performance. Be prepared to refactor data-dependent control flow. | Initial compilation overhead. May require code restructuring for full optimization. [5] [32] |
| Gradient Checkpointing | Drastically reduces memory usage (from $O(n)$ to $O(sqrt(n))$), enabling larger models/batch sizes. | Apply at the block level. Benchmark the number of segments to find the optimal memory/speed trade-off. | Increased computation time due to activation recomputation. [26] [27] |
| Layer Normalization (LN) | Improved training stability, independence from batch size, and avoidance of BN-related issues. | Standard for ConvNeXt; not a tunable hyperparameter but a core architectural choice. | May interact unexpectedly with data preprocessing; requires empirical validation. [17] |
| Custom Normalization | Potential for improved performance by aligning normalization with data characteristics. | Experiment with removing or modifying standard image normalization steps. | Can lead to unexpected behavior if not done carefully; requires rigorous testing. [17] |

In essence, these advanced techniques represent the cutting edge of training optimization. Integrating `torch.compile` is a low-effort, high-reward enhancement that should be a mandatory part of any modern training pipeline. Gradient checkpointing provides a powerful escape hatch when memory constraints arise, but its application should be measured and targeted. Finally, a deeper understanding of the normalization layers within ConvNeXt can reveal subtle but important interactions with data preprocessing that can be exploited to gain a performance advantage. By strategically deploying these tools, you can build a highly optimized and efficient training environment that maximizes the potential of your ConvNeXt models.

# Architectural Variants, Pre-training, and Performance Scaling

The ConvNeXt family of models is not monolithic; it comprises multiple architectural variants and families, each with distinct design philosophies and performance characteristics. Understanding these differences is crucial for selecting an appropriate base model and for designing a systematic investigation into how model properties scale with performance. The evolution from the original ConvNeXt to ConvNeXt V2 represents a significant architectural refinement, while the sheer size of the family, from tiny variants to massive models, provides a rich landscape for studying the scaling laws of deep learning. This section provides a comparative overview of the key architectural variants, explores the profound impact of pre-training strategies, and analyzes the documented scaling relationships between model size, input resolution, and final accuracy on standard benchmarks like ImageNet.

The primary distinction in the ConvNeXt ecosystem is between the original family and the ConvNeXt V2 family [8] [30]. The original ConvNeXt series, including variants like Tiny, Small, Base, Large, and Extra Large, established the modern CNN paradigm by replacing the complex building blocks of older architectures with simpler, more uniform inverted bottlenecks inspired by transformers [10]. The key innovation in ConvNeXt V2 was the introduction of Global Response Normalization (GRN) [30] [34]. GRN is a layer that performs global feature aggregation, divisive normalization, and calibration to increase inter-channel feature diversity and selectivity [34]. It was integrated into the ConvNeXt block, replacing the redundant LayerScale mechanism [34]. The inclusion of GRN was a co-design effort alongside the development of the Fully Convolutional Masked Autoencoder (FCMAE) pre-training framework, which jointly enabled ConvNeXt V2 to achieve state-of-the-art results [30]. For your research, comparing the performance of a model from the original ConvNeXt family against its counterpart in the ConvNeXt V2 family (e.g., ConvNeXt-B vs. ConvNeXt V2-Base) is a scientifically valid and insightful experiment. It would allow you to quantify the contribution of the GRN layer and the FCMAE pre-training methodology to overall performance.

The performance of these models is profoundly influenced by the pre-training strategy. While fine-tuning on large datasets like ImageNet-1K is a common and effective approach, self-supervised pre-training followed by supervised fine-tuning has emerged as a powerful alternative. ConvNeXt V2 models were pre-trained using

the FCMAE framework, where a large fraction of image patches (60%) are masked and the network learns to reconstruct them from the visible ones [34] . This pre-training is followed by a relatively short supervised fine-tuning phase on ImageNet-1K [30] . The results of this approach are compelling. For example, the FCMAE-pretrained ConvNeXt V2-Large achieved an 82.94% top-1 accuracy on ImageNet-1K. When this model was further pre-trained on ImageNet-21k before the final fine-tuning on ImageNet-1K, its accuracy rose to 83.89%, demonstrating the benefit of pre-training on a larger, more diverse dataset [35] . This pattern holds across the entire V2 family, with ImageNet-21k pre-training yielding consistent accuracy gains [35] . For your research, this suggests that a promising direction would be to compare the performance of a model pre-trained on ImageNet-1K with a model pre-trained on ImageNet-21k (if available) or a self-supervised method like FCMAE, followed by fine-tuning. This would provide direct insight into the impact of pre-training on your specific task.

The scaling of ConvNeXt models has been extensively studied, revealing clear trends between model size, computational cost (measured in FLOPs), input resolution, and final accuracy. Generally, larger models achieve higher accuracy. For instance, moving from ConvNeXt-Tiny (28M params) to ConvNeXt-Large (198M params) on ImageNet-1K increases the top-1 accuracy from 82.1% to 84.4% at a 224x224 resolution [6] [7] . Similarly, increasing the input resolution from 224x224 to 384x384 consistently leads to higher accuracy. The ConvNeXt-B model sees its accuracy jump from 83.8% to 85.1%, and the ConvNeXt-XL model improves from 87.0% to 87.8% when the input resolution is increased [7] . This highlights the importance of not only scaling the model's internal dimensions (`dims` list in the architecture definition [10] ) but also increasing the input resolution to unlock the full potential of larger models. The user's plan to create a model with "arbitrary depth that is scaled accordingly" aligns perfectly with these documented scaling principles. A systematic experimental design should explore the joint effects of scaling the `depths` list, the `dims` list, and the input resolution.

The table below provides a comparative overview of the architectural families and their performance characteristics.

| Feature | Original ConvNeXt Family | ConvNeXt V2 Family |
|---|---|---|
| Key Innovation | Modern inverted bottleneck blocks replacing older primitives. | Global Response Normalization (GRN) layer. [10] [30] |
| Normalization | Layer Normalization (LN). | Layer Normalization (LN) with GRN. [17] [30] |
| Pre-training Method | Typically supervised fine-tuning from ImageNet pre-trained weights. | Primarily self-supervised FCMAE pre-training, followed by supervised fine-tuning. [8] [34] |
| Architectural Variants | Tiny, Small, Base, Large, XLarge. | Atto, Femto, Pico, Nano, Tiny, Base, Large, Huge, XXL, 660M. [8] [10] |
| ImageNet-1K Top-1 Acc (224x224) | Tiny: 82.1%, Small: 83.1%, Base: 83.8%, Large: 84.4%. | Atto: 76.7%, Tiny: 78.5%, Base: 81.9%, Large: 82.9%, Huge: 84.9%. [7] [8] [35] |
| ImageNet-21k Pretraining Impact | Available for some variants ('in22ft1k' suffix). | Consistently boosts accuracy when used before ImageNet-1K fine-tuning. [11] [35] |

In conclusion, the choice of architectural variant and pre-training strategy is a critical axis of experimentation. The ConvNeXt V2 family, with its GRN layer and FCMAE pre-training, represents the current state-of-the-art within this architectural lineage. A systematic exploration of these models, coupled with a study of scaling laws related to model size, resolution, and pre-training data, will provide deep insights into the factors that govern performance. Your research can contribute meaningfully by validating these scaling behaviors on your specific image recognition task and potentially discovering novel scaling regimes or architectural modifications.

# Structuring a Rigorous Hyperparameter Sweep for Research

Conducting a hyperparameter sweep is a systematic process that moves beyond random guessing to a structured exploration of the parameter space. A rigorous sweep is designed to identify optimal configurations, understand the sensitivity of the model to different parameters, and ultimately draw reliable conclusions about what drives performance. Based on the comprehensive analysis of the provided research, this final section synthesizes all preceding discussions into a concrete, actionable framework for structuring your hyperparameter sweep. This framework is built around establishing a strong baseline, conducting a systematic search over key parameters, and evaluating the impact of advanced optimizations. The goal is to transform your training process from a black box into a transparent and reproducible scientific experiment.

The first step in any rigorous sweep is to establish a robust baseline configuration. This baseline should be informed by the extensive body of research summarized in the previous sections. A strong starting point would be to use the most commonly adopted and successful combination of hyperparameters found across multiple studies. This includes using theAdamW optimizerwith an initial learning rate of $1\times10^{-4}$, a weight decay of0.01, and acosine annealing learning rate schedulethat includes awarmup period of 5 epochs [3] [14] [15] [20] [21] . For regularization, start with the standardDropPath rate (`drop_path_rate`) of 0.1 [11] [23] . For the training environment, enableAutomatic Mixed Precision (AMP)using `torch.cuda.amp` and implementDistributedDataParallel (DDP)for all multi-GPU runs [1] [25] . Begin with a moderate total batch size, such as 64, and adjust the per-GPU batch size accordingly (e.g., 64 on a single GPU, 32 on two GPUs, etc.). This baseline provides a solid, evidence-based starting point from which you can confidently deviate and measure the impact of specific changes.

Once the baseline is established, the sweep itself can begin. The most effective approach is to adopt a grid search or, more practically, a series of controlled experiments where you vary one primary hyperparameter at a time while keeping all others fixed at their baseline values. This isolates the effect of each parameter. The primary axes for your sweep should be:

1.Learning Rate (LR):This is often the most influential hyperparameter. Conduct a broad sweep of the initial LR, covering at least four orders of magnitude, from $1\times10^{-5}$ to $1\times10^{-2}$. Within this range, pay special attention to the "surge peak" phenomenon described in the literature, where the optimal LR for an adaptive optimizer like AdamW reaches a maximum before decreasing again as the batch size increases [28] . This could lead to the discovery of unconventional but highly effective LR/batch size pairs.

2.Batch Size:As you vary the LR, you must also scale the batch size to maintain a constant data throughput. This is a direct consequence of the LR scaling principle [9] . Test a range of total batch sizes, perhaps from 32 to 1024, and for each size, determine the corresponding per-GPU batch size and number of gradient accumulation steps needed to fit within the 24GB VRAM of the RTX 4090.

3.Weight Decay:Sweep the weight decay parameter across a logarithmic range, such as$1\times10^{-5}$,$1\times10^{-4}$,$1\times10^{-3}$,$1\times10^{-2}$, and $1\times10^{-1}$. This will help you understand the optimal level of regularization for your specific task and model architecture.

4.Model Architecture:Treat the choice of model architecture as a hyperparameter. Compare the original ConvNeXt family against the ConvNeXt V2 family. For each family, test different scales (e.g., Tiny, Base, Large) and resolutions (e.g., 224x224 vs. 384x384) to map out the scaling curve for your task. This involves systematically varying the `depths` and `dims` lists in the model configuration [10].

For each unique configuration in your sweep, you should run the experiment multiple times (e.g., 3-5 repetitions) with different random seeds to account for stochasticity and ensure the results are robust [21]. During each run, monitor key metrics beyond just final accuracy. Track the training and validation loss, convergence speed (number of epochs to reach a target accuracy), and GPU utilization to get a holistic view of the training dynamics. The ultimate goal is not just to find the single best-performing model but to understand the trade-offs involved. For example, a configuration with a slightly lower final accuracy but significantly faster training time and higher GPU utilization might be preferable depending on your research goals and resource constraints.

Finally, the sweep should incorporate evaluations of advanced optimizations. For each baseline model configuration, test the impact of `torch.compile`. Measure the training time with and without compilation to quantify the speedup. Similarly, if memory becomes a constraint, benchmark the effect of gradient checkpointing on both training time and peak memory usage for your deepest model variants [26] [27]. These experiments will provide valuable insights into the performance bottlenecks of your specific training setup and guide future optimization efforts.

To summarize, the following structured approach is recommended for your hyperparameter sweep:

| Step | Action | Key Parameters to Vary | Expected Outcome |
|---|---|---|---|
| 1. Establish Baseline | Implement the most common and successful configuration from literature. | Optimizer (AdamW), LR ($1\times10^{-4}$), Weight Decay (0.01), Scheduler (Cosine Annealing + Warmup). | A stable, reasonably fast training process to serve as a reference point. |
| 2. Isolate Learning Rate | Sweep the initial learning rate over a wide range while keeping other parameters fixed. | Initial LR: $1\times10^{-5}$ to $1\times10^{-2}$. | Identify the optimal LR range and investigate the "surge peak" phenomenon. |
| 3. Scale Batch Size | For each LR, scale the batch size proportionally. | Total Batch Size: 32 to 1024. | Determine the optimal batch size for each LR and assess the impact of LR scaling. |
| 4. Tune Regularization | Sweep the weight decay parameter. | Weight Decay: $1\times10^{-5}$ to $1\times10^{-1}$. | Find the optimal level of regularization to prevent overfitting. |
| 5. Explore Architecture | Compare different ConvNeXt families and scales at different resolutions. | Model Family (ConvNeXt vs. V2), Size (Tiny/Base/Large), Resolution (224x224/384x384). | Map the scaling laws of the architecture on your specific task. |
| 6. Evaluate Advanced Tools | Benchmark `torch.compile` and gradient checkpointing. | `torch.compile`: On/Off. Gradient Checkpointing: Segments=1, 2, 4, ... | Quantify the performance gains from system-level optimizations. |

By following this structured, evidence-based framework, you can move beyond ad-hoc tuning and conduct a rigorous, scientific investigation into the optimal training of your ConvNeXt model. This approach will not only help you find the best possible configuration for your task but will also generate valuable insights into the underlying principles of deep learning optimization.

---

## Reference

1. Get Started with PyTorch 2.4 and CUDA 12.4 on Runpod https://www.runpod.io/articles/guides/pytorch-2-4-cuda-12-4

2. When I was training with multiple GPUs, I kept getting stuck ... https://github.com/ultralytics/ultralytics/issues/11680

3. Optimization and Performance Comparison of AOD-Net ... https://www.mdpi.com/2673-2688/6/8/181

4. Application of ConvNeXt with transfer learning and data ... https://pmc.ncbi.nlm.nih.gov/articles/PMC12136420/

5. Introduction to torch.compile https://docs.pytorch.org/tutorials/intermediate/torch_compile_tutorial.html

6. convnext_large — Torchvision main documentation http://docs.pytorch.org/vision/main/models/generated/torchvision.models.convnext_large.html

7. Code release for ConvNeXt model https://github.com/facebookresearch/ConvNeXt

8. Code release for ConvNeXt V2 model https://github.com/facebookresearch/ConvNeXt-V2

9. Hyperparameters to reproduce ConvNeXt-atto/femto/pico/ ... https://github.com/huggingface/pytorch-image-models/discussions/1671

10. ConvNeXt/models/convnext.py at main · facebookresearch ... https://github.com/facebookresearch/ConvNeXt/blob/main/models/convnext.py

11. ConvNeXt — tfimm 0.1 documentation https://tfimm.readthedocs.io/en/latest/content/convnext.html

12. pprp/timm: PyTorch image models, scripts, pretrained weights https://github.com/pprp/timm

13. Accelerated PyTorch Training on M1 Mac https://news.ycombinator.com/item?id=31424048

14. Instance-level semantic segmentation of nuclei based on ... https://pmc.ncbi.nlm.nih.gov/articles/PMC11804060/

15. NTIRE 2025 Challenge on RAW Image Restoration and ... https://arxiv.org/html/2506.02197v2

16. Brain tumor grade classification using the ConvNext ... https://journals.sagepub.com/doi/10.1177/20552076241284920

17. Computer vision for plant pathology: A review with examples ... https://bsapubs.onlinelibrary.wiley.com/doi/10.1002/aps3.11559

18. A deep learning based smartphone application for early ... https://pmc.ncbi.nlm.nih.gov/articles/PMC11685909/

19. The State of Machine Learning Competitions https://mlcontests.com/state-of-machine-learning-competitions-2024/

20. Instance-level semantic segmentation of nuclei based on ... https://bmcbioinformatics.biomedcentral.com/articles/10.1186/s12859-025-06066-8

21. A Critical Assessment of Generative Models for Synthetic Data ... https://pmc.ncbi.nlm.nih.gov/articles/PMC10741143/

22. GPUs for Large Language Models: Kernels, Triton, Memory ... https://medium.com/@hexiangnan/gpus-for-large-language-models-kernels-triton-memory-coalescing-and-the-execution-hierarchy-7aaa32dac5ae

23. InceptionNeXt: When Inception Meets ConvNeXt (CVPR ... https://github.com/sail-sg/inceptionnext

24. GPU Performance Deep Learning Benchmarks | Exxact Blog https://www.exxactcorp.com/blog/benchmarks/gpu-performance-deep-learning-benchmarks

25. Benchmarking Advanced Multi – GPU Training Strategies https://medium.com/@savyasachi.thati/benchmarking-advanced-multi-gpu-training-strategies-20c9675003db

26. Saving memory using gradient-checkpointing https://github.com/cybertronai/gradient-checkpointing

27. Gradient checkpointing and its effect on memory and runtime https://discuss.pytorch.org/t/gradient-checkpointing-and-its-effect-on-memory-and-runtime/198437

28. Surge Phenomenon in Optimal Learning Rate and Batch … https://arxiv.org/html/2405.14578v2

29. Explaining neural scaling laws - PMC https://pmc.ncbi.nlm.nih.gov/articles/PMC11228526/

30. Review — ConvNeXt V2: Co-designing and Scaling ConvNets … https://sh-tsang.medium.com/review-convnext-v2-co-designing-and-scaling-convnets-with-masked-autoencoders-4346eadb5405

31. YOLO-PLNet: a lightweight real-time detection model for … https://pmc.ncbi.nlm.nih.gov/articles/PMC12669193/

32. Impact of ML optimization tactics on greener pre-trained ML … https://link.springer.com/article/10.1007/s00607-025-01437-8

33. LLM as a Neural Architect: Controlled Generation of Image … https://www.researchgate.net/publication/398520845_LLM_as_a_Neural_Architect_Controlled_Generation_of_Image_Captioning_Models_Under_Strict_API_Contracts

34. Papers Explained 94: ConvNeXt V2 https://medium.com/thedeephub/papers-explained-94-convnext-v2-2ecdabf2081c

35. ConvNeXt V2 — MMPretrain 1.2.0 documentation https://mmpretrain.readthedocs.io/en/latest/papers/convnext_v2.html