

## The V-Shape Breakthrough: How AI Researchers Cut Model Training Memory in Half—And Made It Faster

Training massive AI models like GPT-4 is an exercise in extremes. The process is a constant, expensive battle against computational limits, forcing engineers into a difficult trade-off: you can optimize for speed, or you can optimize for memory, but you can rarely have both. Pushing for faster training usually demands more memory, while cutting memory usage typically slows everything down. This fundamental tension has defined the landscape of large-scale AI for years.

But what if this trade-off isn't as fundamental as we thought? A team of researchers has introduced a new framework that reframes this problem in a surprisingly elegant way. Instead of making incremental tweaks to existing methods, they went back to the first principles of how training tasks are scheduled. By focusing on a simple, repeating "building block" at the heart of the process, they uncovered a new path to efficiency.

The results are deeply counter-intuitive and incredibly promising. This new approach leads to training schedules that can cut peak memory usage in half while simultaneously *increasing* training speed. It also reveals a method to achieve maximum theoretical throughput eliminating wasteful GPU idle time without paying the expected memory penalty. This isn't just a minor improvement; it's a new blueprint for efficiency.

### 1. The Big Insight: It's All About Balancing Task 'Lifespans'

The paper's core conceptual breakthrough is to view a complex pipeline training schedule not as a chaotic sequence of operations, but as the simple repetition of a core "building block." This block contains the set of forward and backward passes for a single chunk of data. How this block is designed and repeated determines the efficiency of the entire system.

The researchers discovered a critical link: the duration, or "lifespan," of this building block directly dictates the peak memory required during training. The lifespan is the time from when a piece of data's activation memory is created (at the start of a forward pass) until it's no longer needed (after the backward pass). As the paper concisely states lifespan decides the activation memory.

This insight reveals why most existing methods are "memory inefficient." They create an imbalanced workload where the first device in the training pipeline is stuck with tasks that have the longest lifespans. This creates a severe memory bottleneck on that single device, while the memory on all other devices is underutilized.

The elegant solution is a new design pattern called the "V-Shape" building block. Its power comes from workload balancing. Stages with long activation lifespans are paired on the same device with stages that have short lifespans, ensuring that the total memory load is balanced evenly across all devices. This is achieved by partitioning the model into twice the number of stages as there are devices and then reversing the device placement for the second half of the stages. For example, on 4 devices, stages 1-4 are placed on devices 1-4, while stages 5-8 are placed on devices 4-1 in reverse order creating the characteristic "V" shape.

## 2. The Counter-Intuitive Win: Cut Memory by 50% and Actually Go *Faster*

Applying this V-Shape principle led to the creation of the "V-Half" schedule, a new method with a remarkable outcome. V-Half reduces the peak activation memory to approximately half of what the widely-used "1F1B" schedule requires—a massive saving in the memory-constrained world of large model training. (The paper notes the exact peak memory is  $\lceil \frac{d+1}{2} \rceil M/d$ , which approaches  $1/2$  as the number of devices  $d$  increases).

Conventionally, engineers would resort to a technique called activation rematerialization to free up memory, which involves re-computing information on the fly. This saves memory but burns valuable GPU cycles, significantly slowing down training. V-Half provides a far more elegant solution, saving memory without this computational penalty. V-Half is one of several schedules enabled by the V-Shape framework. The researchers also developed **V-Min**, which can cut activation memory down to just *one-third* of 1F1B, albeit with a more significant impact on throughput.

The most surprising result, however, is that this huge memory saving doesn't come with a performance penalty it can make training *faster*. This counter-intuitive win is possible because the memory freed up by the V-Half schedule can be repurposed. With more available memory, the system can process larger "microbatches" of data at once.

Processing larger microbatches increases the system's "arithmetic intensity." In simple terms, this means the GPUs spend more of their time performing powerful mathematical calculations (the "arithmetic") and less time waiting for smaller chunks of data to be moved around (the "intensity"). This shift directly translates into a significant throughput boost. The paper's experiments confirm this, showing that on bigger models where memory is extremely limited, the V-Half schedule can surpass other methods for this exact reason.

## 3. The 'Free Lunch': Achieve Peak Speed for the Same Memory Price

Beyond just saving memory, the V-Shape framework also unlocked a way to maximize speed. The researchers developed another schedule called "V-ZB" (V-Zero Bubble), which tackles one of the biggest sources of inefficiency in pipeline parallelism: "pipeline bubbles." These bubbles are periods of inactivity where expensive GPUs sit idle, waiting for the next task. They are pure waste, costing time and money.

The breakthrough of the V-ZB schedule is that it achieves "almost zero pipeline bubbles while maintaining the same activation memory as 1F1B." This is a massive achievement because it effectively breaks the traditional speed-memory trade-off. It delivers a huge performance gain pushing throughput to its theoretical maximum without the expected cost of higher memory usage.

This is what makes V-ZB so important. It offers something akin to a free lunch for AI training. As the paper's authors put it, this is "almost a pure gain to the existing methods." The paper notes that this massive gain comes at the minor cost of doubled communication between pipeline

stages a trade-off the authors describe as "relatively small and can be neglected" in most practical scenarios.

### Conclusion: A New Blueprint for Efficiency

This research offers more than just a few new scheduling techniques; it provides a powerful new framework for thinking about and designing efficient AI training systems. By deconstructing the problem into its fundamental "building blocks" and focusing on the concept of "lifespan," the researchers have created a systematic way to control the trade-offs between memory and speed.

The V-Half and V-ZB schedules are potent demonstrations of this framework's power, but the underlying methodology is the real prize. The researchers even developed an **adaptive scheduler** that can search for optimal configurations, giving engineers a set of tools to create bespoke pipelines that are perfectly tuned for any memory limit, model, or hardware configuration. This work shows how rethinking a problem from its most basic structure can unlock massive gains. What other accepted trade-offs in AI are just waiting for a similar 'V-Shape' insight to break them wide open?

### **Briefing on Pipeline Parallelism with Controllable Memory**

This document synthesizes findings on a novel framework for designing and analyzing pipeline parallelism (PP) schedules for training large-scale models. The central problem addressed is that existing PP methods, while crucial for scaling model training, suffer from significant inefficiencies, namely **pipeline bubbles** (compute idle time) and **high activation memory consumption**.

The core contribution is a systematic framework that deconstructs any pipeline schedule into a repeating "**building block**." This framework reveals a fundamental insight: the "**lifespan**" of this building block—the time from its first forward pass to its last backward pass directly determines the schedule's peak activation memory. Analysis using this framework shows that nearly all existing pipeline schedules are memory-inefficient due to imbalanced memory distribution across devices and redundant data dependencies.

To resolve this, a new family of "**V-Shape**" **building blocks** is introduced. These schedules strategically partition the model and place stages in a V-shaped order across devices, co-locating long-lifespan and short-lifespan stages to balance peak memory usage. This design allows for controllable memory-throughput trade-offs. The three primary variants are:

- **V-Min:** Reduces peak activation memory to approximately **1/3** of the standard 1F1B schedule.
- **V-Half:** Reduces peak activation memory to approximately **1/2** of 1F1B while improving throughput.
- **V-ZB:** Achieves **almost zero pipeline bubbles**, maximizing throughput with memory usage comparable to 1F1B.

Experimental evaluations on large language models (up to 98.5B parameters) demonstrate significant performance gains. In pure pipeline parallelism settings, the proposed methods outperform the 1F1B baseline by **7% to 55%** in throughput. In practical, hybrid parallelism scenarios involving grid search over hyperparameters, the V-Shape schedules deliver a **16% throughput improvement** over 1F1B, establishing a new Pareto frontier for memory efficiency and computational throughput in large model training.

## 1. The Challenge of Pipeline Parallelism in Large Model Training

Pipeline parallelism (PP) is a foundational model parallelism strategy for distributed training of large language models. It addresses memory constraints by partitioning a model's parameters across a set of devices, which process microbatches in a streaming fashion. PP is particularly advantageous in environments where communication bandwidth is a bottleneck, as its communication costs are relatively low compared to alternatives like Tensor Parallelism (TP). Despite its wide adoption, PP suffers from two prominent disadvantages:

- **Pipeline Bubbles:** These are periods of inactivity on compute devices that occur during the "warm-up" and "cool-down" phases of the pipeline. Eliminating these bubbles is a primary goal for improving computational efficiency. While asynchronous PP methods can theoretically eliminate bubbles, they do so at the cost of compromising optimization semantics, which can harm model convergence.
- **Large Activation Memory:** Synchronous PP methods must store the intermediate activations from the forward pass of multiple microbatches to be used during the backward pass. This creates a significant memory footprint. Methods like GPipe increase the number of microbatches to reduce the bubble rate, but this exacerbates activation memory consumption. The widely-used 1F1B schedule avoids this memory growth but retains the same bubble rate as GPipe. Existing solutions like GEMS, Chimera, BPipe, and Zero Bubble have attempted to address these issues, but often introduce trade-offs such as large bubble rates, doubled parameter memory, or high communication overhead. The presented work argues that these methods lack a systematic methodology, resulting in memory-inefficient designs.

## 2. A Deconstructive Framework for Pipeline Schedules

A novel four-step framework is proposed to formalize the design and analysis of pipeline schedules. This approach posits that most existing schedules can be understood as variations of a repeating core pattern.

1. **Building Block:** The process begins by defining the sequence of passes (Forward, Backward) for a single microbatch. This sequence is the fundamental "building block" of the entire schedule.
2. **Repeating:** The building block is repeated for subsequent microbatches, woven together to form the pipeline. A legitimate building block must be repeatable without collisions, meaning passes from different microbatches do not overlap on the same device at the same time.

**3. Squeezing:** After repeating the blocks, the resulting schedule may contain redundant idle time (bubbles). A "squeezing" step compacts the schedule by removing these gaps without altering the relative order of passes.

**4. Reordering (Optional):** The warm-up and cool-down phases of the pipeline often have underutilized memory. Passes in these phases can be reordered to further improve computation throughput without affecting the peak memory, which typically occurs in the stable phase.

### The Core Insight: Lifespan Determines Peak Memory

The framework's primary analytical power comes from a direct link identified between a building block's properties and the pipeline's peak memory consumption. Two quantities are crucial:

- **Lifespan (1):** The amount of time between the start of a stage's forward pass (F) and the end of its corresponding backward pass (B/W). Activation memory for a microbatch is allocated at the start of F and held until the end of this lifespan.
- **Repeating Interval (T):** The time duration of the building block on a single stage. In an efficient pipeline, this equals the units of computation per stage within the block.  
The peak memory for a device can be calculated with a simple formula, which sums the contributions from all stages ( $s$ ) allocated to that device ( $i$ ):

$$\text{Peak Memory of Device } i \leq \sum_{s \in S_i} \lceil 1s / T \rceil ms$$

This equation reveals the central thesis: **lifespan decides activation memory**. By designing building blocks with specific lifespans, it becomes possible to create pipeline schedules with controllable and predictable peak memory usage.

### 3. Uncovering Inefficiencies in Existing Pipeline Schedules

Applying this framework reveals that, to the authors' knowledge, all existing pipeline schedules are memory inefficient for two primary reasons:

1. **Redundant Dependency Chain:** Most schedules treat the backward pass as a monolithic operation. By splitting the backward pass into two distinct computations one for activation gradients (**B**) and one for weight gradients (**W**), as introduced in the Zero Bubble paper the lifespan can be shortened, as activation memory can be freed after the **B** pass is complete.
2. **Imbalanced Memory Usage:** In typical schedules like 1F1B, the lifespans of different stages are highly heterogeneous. The first stage has the longest lifespan, creating a memory bottleneck on the first device while subsequent devices have underutilized memory.

### 4. V-Shape Building Blocks: A Memory-Efficient Solution

To address the issue of imbalanced memory, a family of novel, memory-efficient building blocks called **V-Shape schedules** is proposed.

## The V-Shape Concept

The core idea is to balance the sum of lifespans on each device. This is achieved through a specific model partitioning and device placement strategy:

- The model is partitioned into  $2d$  stages, where  $d$  is the number of devices.
- The first  $d$  stages are placed sequentially on devices 1 to  $d$ .
- The second  $d$  stages are placed in reverse order on devices  $d$  down to 1.

This "V-Shape" placement ensures that each device hosts one stage with a long lifespan (from the first half) and one with a short lifespan (from the second half). This collocation balances the total lifespan per device, eliminating the memory bottleneck seen in traditional linear placements.

## Controllable Memory and Schedule Variants

The V-Shape design enables direct control over peak memory by adjusting the time offsets ( $\delta_0, \delta_1$ ) between adjacent Forward and Backward passes across devices. This leads to three distinct, representative schedules:

Building Block	$\delta_0$	$\delta_1$	Peak Activation Memory	Pipeline Bubbles	Description
<b>1F1B (Baseline)</b>	2	4	$M$	$\sim 6d$	Standard memory-imbalanced schedule.
<b>V-Min</b>	1	1	$\approx M/3$	$\approx 4d$	Minimizes memory but can have repeating bubbles if F/B/W times differ.
<b>V-Half</b>	2	1	$\approx M/2$	$\approx 3d$	A balanced option with significantly less memory and fewer bubbles than 1F1B.
<b>V-ZB</b>	4	2	$M$	0	Eliminates bubbles for maximum throughput at memory cost similar to 1F1B.

*Note:  $M$  is the total activation memory of the model, and  $d$  is the number of devices.*

A key limitation of **V-Min** is its susceptibility to "repeating bubbles" when the run times of F, B, and W passes are not equal, causing bubble rates to grow with the number of microbatches. **V-Half** is more robust to such variations and tessellates well in most empirical cases.

## 5. Experimental Validation and Performance Gains

The proposed V-Shape schedules were implemented in the Megatron-LM framework and evaluated on up to 40 NVIDIA A100 GPUs using GPT-3-like models ranging from 9.6B to 98.5B parameters.

## Pure Pipeline Parallelism Results

- **Throughput and Memory:** In pure PP settings, V-ZB consistently achieved the highest throughput (measured in MFU), outperforming all other methods. V-Min and V-Half successfully reduced activation memory to approximately 1/3 and 1/2 of 1F1B, respectively, while V-Min maintained comparable throughput to 1F1B for a smaller number of microbatches. Together, the V-Shape schedules establish a new **Pareto frontier** of throughput versus memory consumption.
  - **The Advantage of Memory Savings:** The memory freed up by V-Min and V-Half allows for using larger microbatch sizes, which increases arithmetic intensity. Experiments show that on larger models with high memory pressure, **V-Half with a doubled microbatch size can surpass the throughput of V-ZB**, demonstrating a practical benefit of memory efficiency.
- Performance in Practical Hybrid Parallelism Scenarios

To assess performance in realistic training scenarios, a grid search was performed over hybrid parallelism strategies, combining PP with Tensor Parallelism (TP) and Data Parallelism (DP) for a 98.5B model. The results underscore the adaptive strengths of the V-Shape family.

Model & Setup	PP Method	Best MFU (%)	Best Parameters (DP, TP, PP, mbs)
<b>98.5B Model</b>	1F1B	54.77	(2, 4, 5, 2)
Seq Length 1024	V-Half	57.83	(2, 1, 20, 1)
(Lower Memory Pressure)	<b>V-ZB</b>	<b>63.31</b>	(1, 1, 40, 1)
<b>98.5B Model</b>	1F1B	62.95	(2, 4, 5, 1)
Seq Length 3072	V-ZB	62.56	(1, 4, 10, 1)
(Higher Memory Pressure)	<b>V-Half</b>	<b>66.34</b>	(1, 2, 20, 1)
<b>98.5B Model</b>	1F1B	OOM	-
Seq Length 16384	V-ZB	OOM	-
(Extreme Memory Pressure)	<b>V-Half</b>	<b>57.85</b>	(1, 8, 5, 1)

The key takeaways are:

- When memory is less constrained (shorter sequence length), **V-ZB** excels due to its bubble elimination.
- As memory pressure increases (longer sequence length), **V-Half's** memory efficiency becomes critical. It enables configurations (e.g., lower TP degree) that are out-of-memory (OOM) for 1F1B and V-ZB, ultimately achieving the highest throughput.

## 6. Additional Contributions and Future Directions

Beyond the primary V-Shape schedules, the work provides several other valuable contributions:

- **Adaptive Scheduler:** A practical search-based algorithm is proposed to find a pipeline schedule that minimizes bubbles for any given activation memory limit, offering fine-grained control.
- **Generalizability of the Framework:** The building block methodology is shown to be general, leading to the design of other novel schedules, such as **1F1B-V**, which achieves 2/3 of 1F1B's activation memory without requiring the B-W split.
- **Future Work:** The primary limitation identified is the "repeating bubble" issue in V-Min under non-uniform task runtimes. Future research aims to resolve this, potentially by exploring continuous offsets or finer-grained time discretization to further reduce memory consumption without sacrificing scalability.

### **Study Guide for Pipeline Parallelism with Controllable Memory**

1. What are the two prominent disadvantages of traditional Pipeline Parallelism (PP) that the research aims to address?
2. Describe the four-step framework proposed for designing pipeline schedules.
3. What is the core insight presented in the paper that links a pipeline's building block to its peak activation memory?

4. Explain the purpose of splitting the backward pass into "B" (backward for activations) and "W" (backward for weights).
5. What is a "V-Shape" building block and what is its primary design principle for achieving memory efficiency?
6. Briefly compare the primary trade-offs of the three proposed V-Shape schedules: V-Min, V-Half, and V-ZB.
7. Under what conditions does the V-Min schedule suffer from a "repeating bubble," and why is V-Half more robust to this issue?
8. How can memory-saving schedules like V-Half sometimes achieve higher throughput than a zero-bubble schedule like V-ZB?
9. What is the key difference in how Pipeline Parallelism (PP) and Tensor Parallelism (TP) handle activation memory?
10. Besides the V-Shape family, name one other memory-efficient building block mentioned in the paper and its key feature.

### Essay Questions

1. Discuss the proposed four-step framework for designing pipeline schedules. How does this framework reveal the memory inefficiency of prior methods like 1F1B and enable the systematic design of new, memory-efficient schedules?
2. Analyze the design principles of the "V-Shape" building block. Explain how collocating stages with long and short lifespans and controlling offsets ( $\delta$ ) between passes allows for a family of schedules with controllable, balanced memory.
3. Compare and contrast the three primary V-Shape schedules (V-Min, V-Half, and V-ZB) with the 1F1B baseline. Discuss their respective trade-offs in terms of peak activation memory, pipeline bubbles, and throughput, referencing both theoretical analysis and experimental results from the paper.
4. Explain the concept of "repeating bubbles" as it applies to the V-Min schedule. Why does this phenomenon occur in real-world scenarios where F, B, and W passes have different run times, and how do the V-Half and V-ZB schedules mitigate or avoid this issue?
5. The paper argues that its methods push the preference from Tensor Parallelism (TP) to Pipeline Parallelism (PP). Evaluate this claim by discussing the trade-offs between TP and PP in terms of activation memory, communication costs, and single-pass FLOPS utilization (MFU), as detailed in the source context.

### Answer Key

1. Traditional Pipeline Parallelism (PP) suffers from two prominent disadvantages: pipeline bubbles, which are idle times that reduce computational efficiency, and large activation memory consumption. Existing methods often improve one at the cost of the other, for instance, by increasing microbatches to reduce bubbles, which in turn increases memory usage.
2. The proposed four-step framework consists of: 1) defining a **Building Block** for a single microbatch; 2) **Repeating** the block for more microbatches without collision; 3) **Squeezing** the schedule to remove redundant bubbles; and 4) optionally **Reordering** passes in the warm-up and cool-down phases to improve throughput.
3. The core insight is that the peak activation memory of a pipeline schedule is directly decided by the **lifespan** of its building block. The lifespan is the time from the start of a stage's forward pass to the end of its backward pass, and peak memory can be calculated from the ratio of the lifespan to the repeating interval.
4. Splitting the backward pass into 'B' (activation gradients) and 'W' (weight gradients) removes a redundant dependency chain. This allows for a shorter lifespan for the activation memory, as it only needs to be retained until consumed by the 'B' pass, leading to a smaller memory footprint.
5. A "V-Shape" building block partitions the model into twice the number of stages as devices and places the second half of stages in reverse order on the devices. This collocates stages with long lifespans with those of short lifespans, balancing the peak memory load across all devices and avoiding a bottleneck on the first device.
6. V-Min offers the minimum memory consumption (about 1/3 of 1F1B) with comparable throughput but can suffer from repeating bubbles. V-Half reduces memory to about 1/2 of 1F1B with higher throughput. V-ZB achieves almost zero bubbles and the highest throughput while maintaining the same activation memory as 1F1B.
7. V-Min suffers from a "repeating bubble" in real-world scenarios when the run times for F, B, and W passes differ, causing bubbles to grow as the number of microbatches increases. V-Half has looser dependencies and generates patterns that tessellate well in most empirical cases, making its throughput robust to variations in run times.
8. By saving activation memory, V-Half allows for the use of a larger microbatch size, which increases arithmetic intensity. On larger models where memory pressure is high, this gain in arithmetic intensity can compensate for V-Half's inherent bubbles and lead to a higher overall throughput than V-ZB, which may be constrained to a smaller microbatch size.
9. Pipeline Parallelism (PP) typically needs to cache activations for a number of microbatches proportional to the number of stages, resulting in a total activation memory demand similar to an unpartitioned model. In contrast, Tensor Parallelism (TP), when used with sequence parallelism, partitions most activation memory equally among the TP devices.

10. One alternative building block is **1F1B-V**, which applies the V-Shape device placement principle to the 1F1B schedule without splitting the backward pass, reducing peak memory to asymptotically 2/3 of 1F1B. Another is a variation of interleaved 1F1B that achieves the same bubble rate but with lower memory consumption.

## Glossary

Term	Definition
<b>1F1B</b>	A pipeline schedule that avoids activation memory growth with respect to the number of microbatches by staggering forward and backward passes. It serves as a primary baseline for comparison.
<b>Activation Memory</b>	The memory required to store the intermediate results (activations) of the forward passes, which are needed for the backward passes. High activation memory is a major constraint in training large models.
<b>B Pass</b>	The part of the backward pass responsible for computing activation gradients. Splitting this from the W pass can reduce the lifespan of activation memory.
<b>Building Block</b>	The foundational layout of computation passes (F, B, W) for a single microbatch. The paper's framework proposes that entire pipeline schedules can be constructed by repeating and arranging these blocks.
<b>F Pass</b>	The forward pass of a model stage for a single microbatch.
<b>FLOPS Utilization (MFU)</b>	A measure of a training system's throughput, representing the percentage of theoretical peak Floating-Point Operations Per Second (FLOPS) that is actually achieved.
<b>GPipe</b>	An early pipeline parallelism work that reduces the bubble rate by increasing the number of microbatches, but at the cost of significantly more activation memory.
<b>Lifespan</b>	The amount of time between the start of a stage's forward (F) pass and the end of its backward (B or W) pass for a single microbatch. The paper identifies this as the key determinant of peak activation memory.
<b>Microbatch</b>	A small subdivision of a training data minibatch. Pipeline parallelism processes the model in a streaming fashion, with different stages operating on different microbatches simultaneously.
<b>Model Parallelism (MP)</b>	A strategy to partition a model's parameters across a set of devices to address memory constraints when a model is too large for a single device. PP and TP are types of MP.
<b>Pipeline Bubbles</b>	Idle time in the pipeline schedule where devices are not performing useful computation. These bubbles reduce overall training efficiency and throughput.
<b>Pipeline Parallelism (PP)</b>	A model parallelism strategy that splits a model into several sequential stages, which are then processed by different devices in a streaming or pipelined manner.

<b>Rematerialization</b>	A technique to save activation memory by recomputing activations during the backward pass instead of storing them. This saves memory at the cost of repeated computation, which decreases throughput.
<b>Repeating Interval (T)</b>	The time duration of the repeating pattern in the stable phase of a pipeline. In an efficient pipeline, this interval is equal to the number of computation units in each stage of the building block.
<b>Tensor Parallelism (TP)</b>	A model parallelism strategy that partitions weight parameters (tensors) into several devices and performs matrix multiplication separately. It incurs high communication volume but can reduce activation memory.
<b>V-Half</b>	A V-Shape schedule that consumes about half the activation memory of 1F1B with higher throughput. It sits between V-Min and V-ZB in terms of memory savings and bubble reduction.
<b>V-Min</b>	A V-Shape schedule designed for minimum memory consumption (asymptotically 1/3 of 1F1B) with comparable throughput to 1F1B. It is susceptible to "repeating bubbles" if F/B/W times are unequal.
<b>V-Shape Building Block</b>	A novel building block design that partitions the model into twice the number of stages as devices (2d) and places the second half of stages in reverse order. This balances the sum of lifespans on each device, leading to higher memory efficiency.
<b>V-ZB</b>	A V-Shape schedule that eliminates almost all pipeline bubbles, achieving maximum throughput. It requires the same amount of activation memory as the 1F1B baseline.
<b>W Pass</b>	The part of the backward pass responsible for computing weight gradients.