# Deep Learning Frameworks
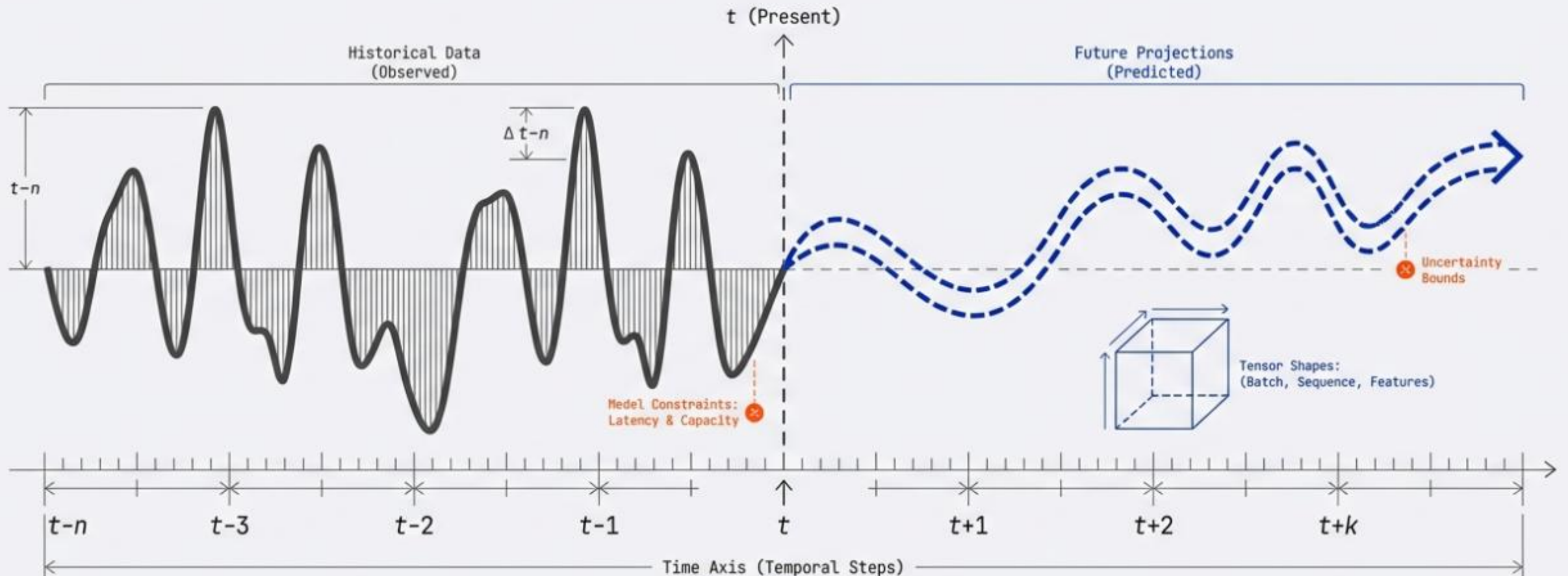
Time Series Forecasting using GRU in Tensorflow, Quantization in Pytorch

# Time Series Forecasting
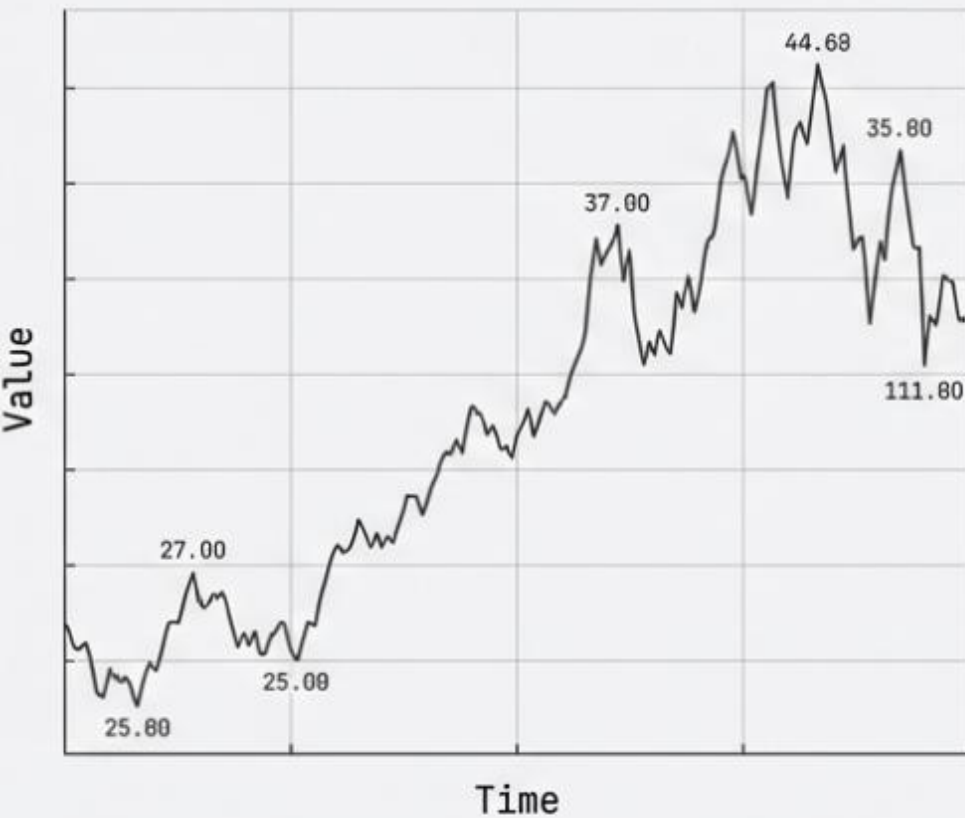


From Raw Data to Future Insight: A Structural Guide for Model Design
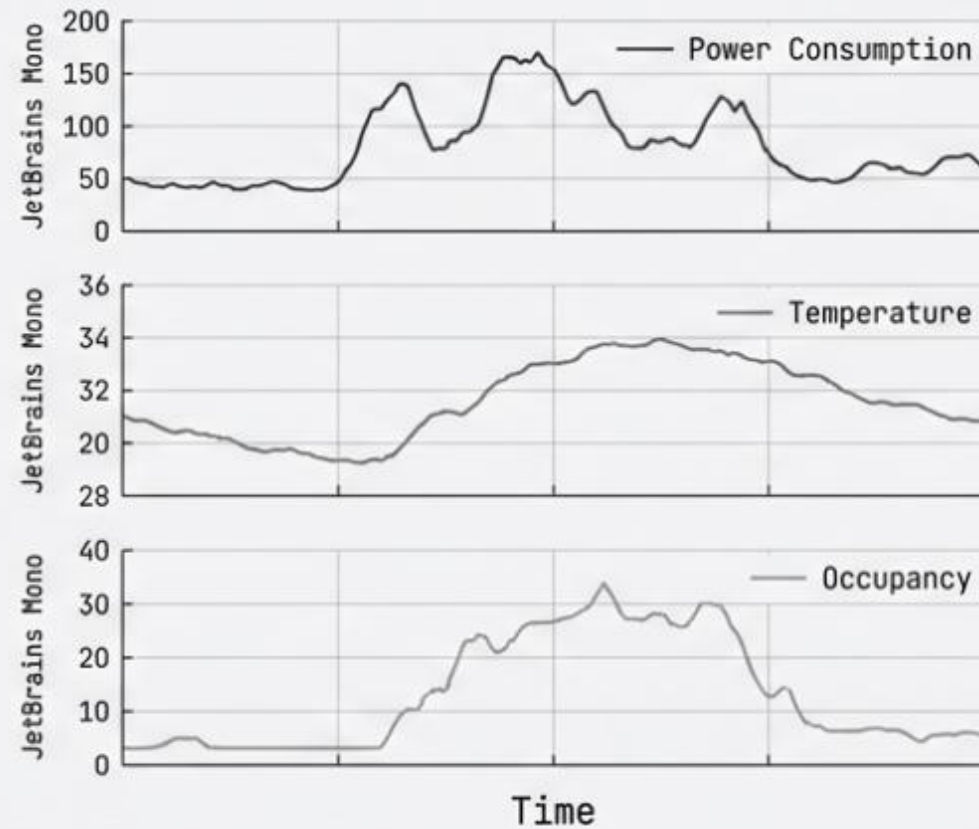
# Univariate and Multi variate



**Univariate Data**
One variable (e.g., Stock Price)

44.68
37.00
35.80
111.80
27.00
25.00
25.80

Value
Time

**Multivariate Data**
Multiple variables (e.g., Energy System)

JetBrains Mono — Power Consumption
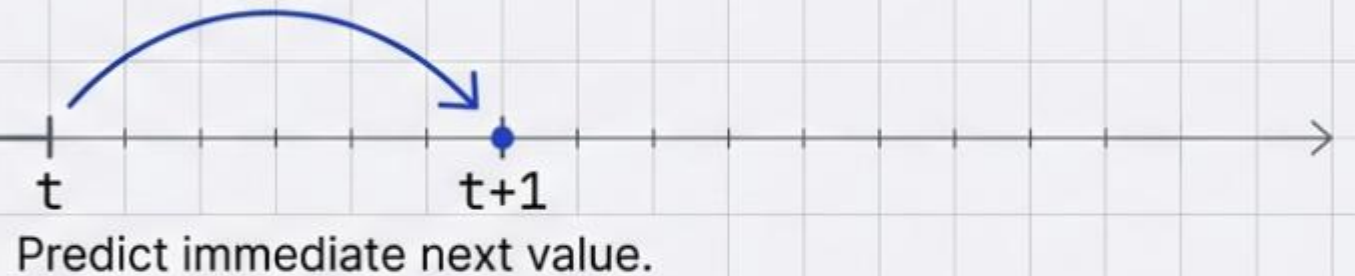JetBrains Mono — Temperature
JetBrains Mono — Occupancy
Time

**Design Decision: Frequency**

Every series requires a fixed sampling frequency (Hourly vs. Daily).
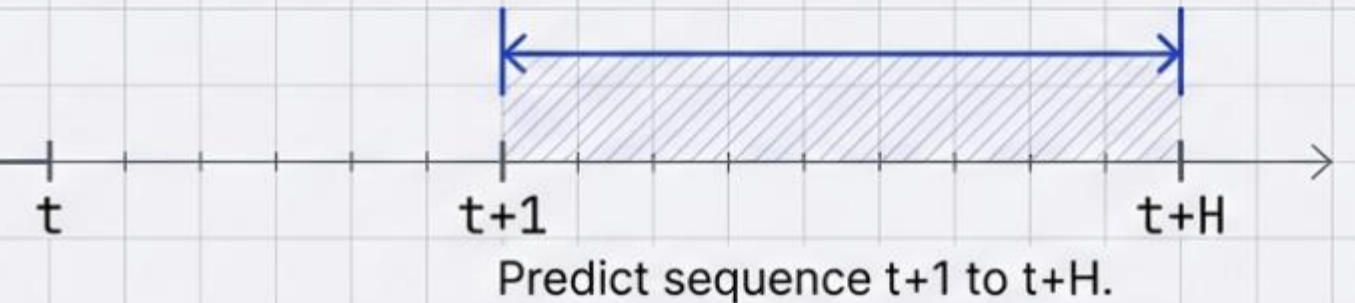This choice dictates visible patterns and input window size.

# Forecast Horizon



The Horizon (H) is the set of future time points the model must generate.

**One-Step Ahead**

t    t+1

Predict immediate next value.

**Multi-Step Ahead**

t    t+1    t+H

Predict sequence t+1 to t+H.

Prediction = $\left[\begin{array}{c} \textbf{Autoregressive} \\ \text{Past values of the Target itself.} \\ \text{(e.g., Past Temp -> Future Temp)} \end{array}\right]$ + $\left[\begin{array}{c} \textbf{Exogenous Variables} \\ \text{External factors influencing the target.} \\ \text{(e.g., Humidity, Holidays, Promotions)} \end{array}\right]$

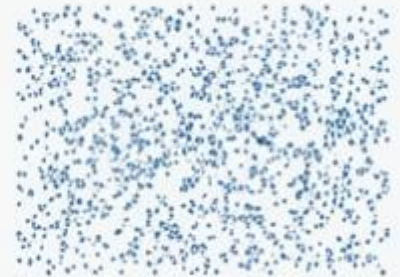# y(t) = Trend + Seasonality + Noise
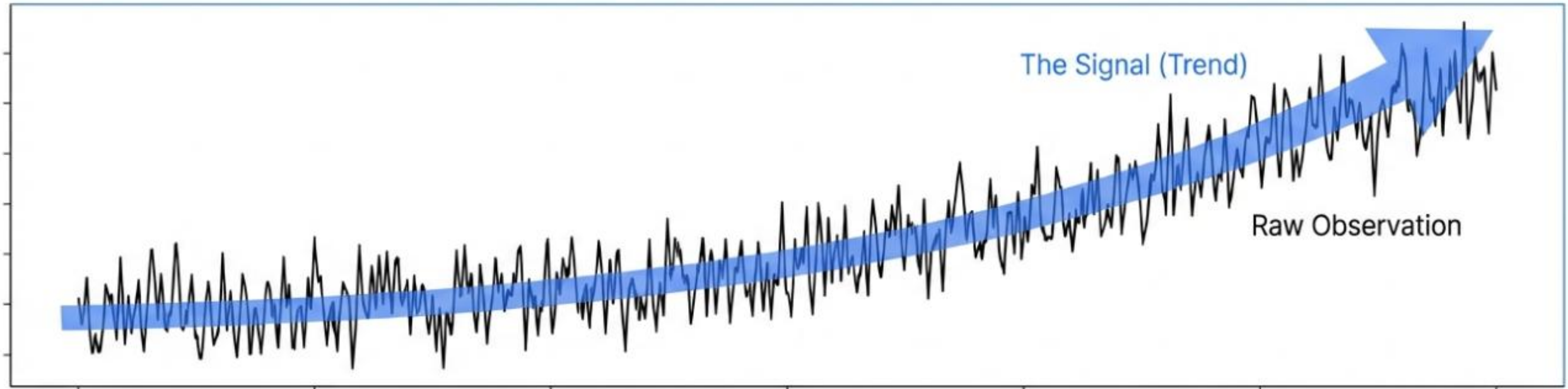
Long-term
direction

Repeating pattern
with a fixed period

Random fluctuations
/ Residuals

# Trend: The long-term Direction



## Definition

The slow, long-term movement of the series.

## Real-World Examples

- **Commercial**: Sales increasing over months as product gains popularity.
- **Physical**: Sensor readings drifting due to mechanical wear.
- **Environmental**: Gradual warming or cooling over seasons.

## Forecasting Intuition

If a series has a strong trend, the model's primary job is learning "direction". While classical models often require manual detrending, deep models can often learn the trend directly if the dataset is sufficiently large.
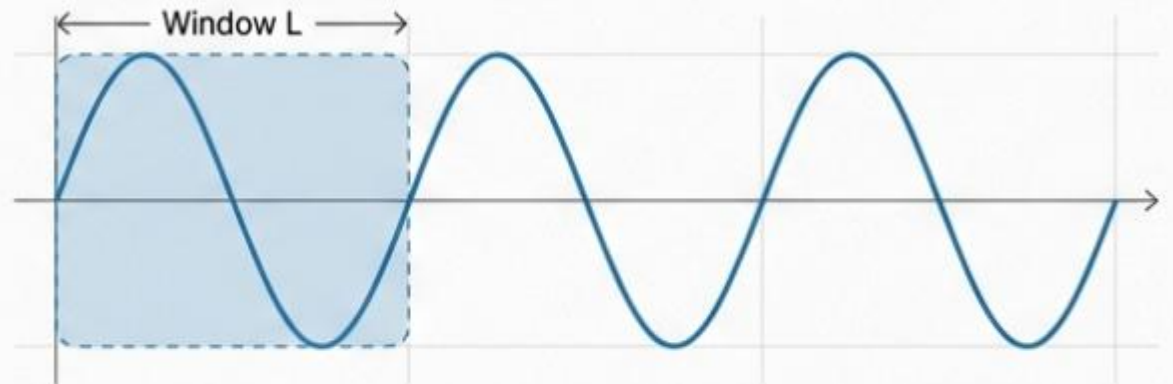
# Seasonality

## Concept

Definition: A repeating pattern with a fixed period.

- Hourly electricity usage → Daily cycle (24 hours)
- Website traffic → Weekly cycle (7 days)
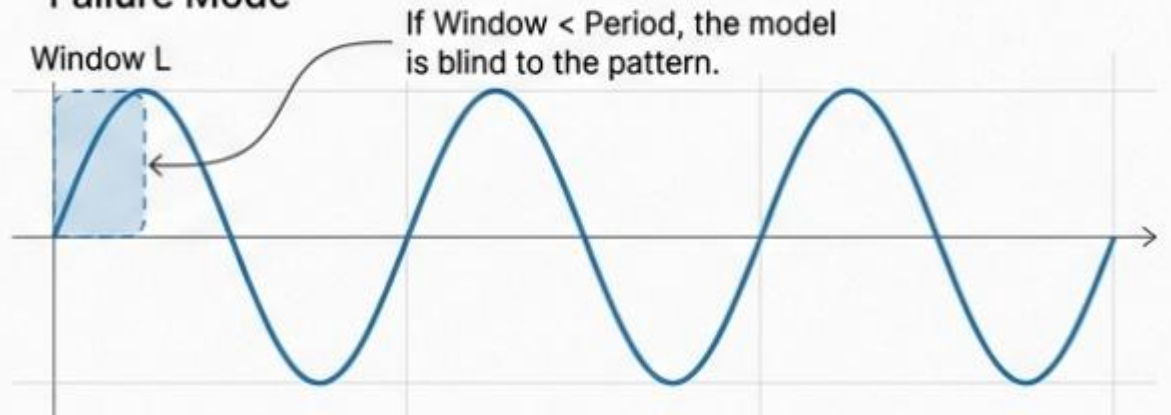- Retail sales → Yearly cycle (festival seasons)

**Forecasting Intuition:** Seasonality means "the past repeats." Your lookback window "L" must be at least one full season.

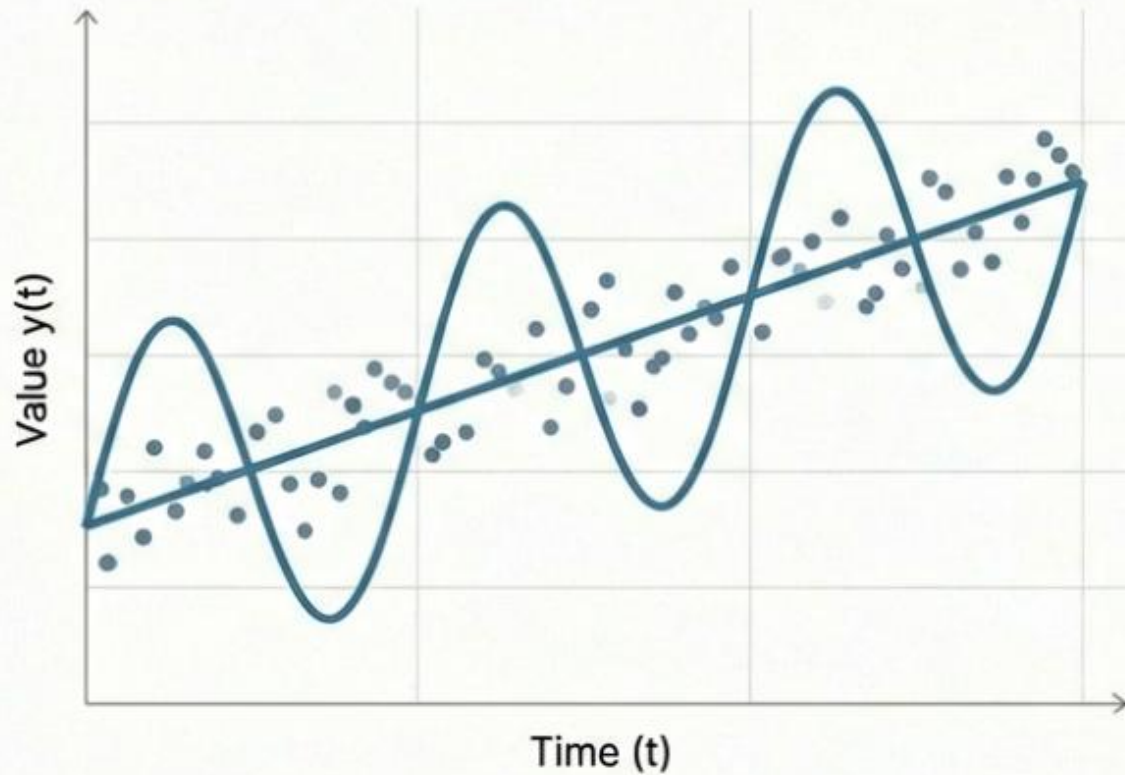## The Architectural Constraint

### Window Rule
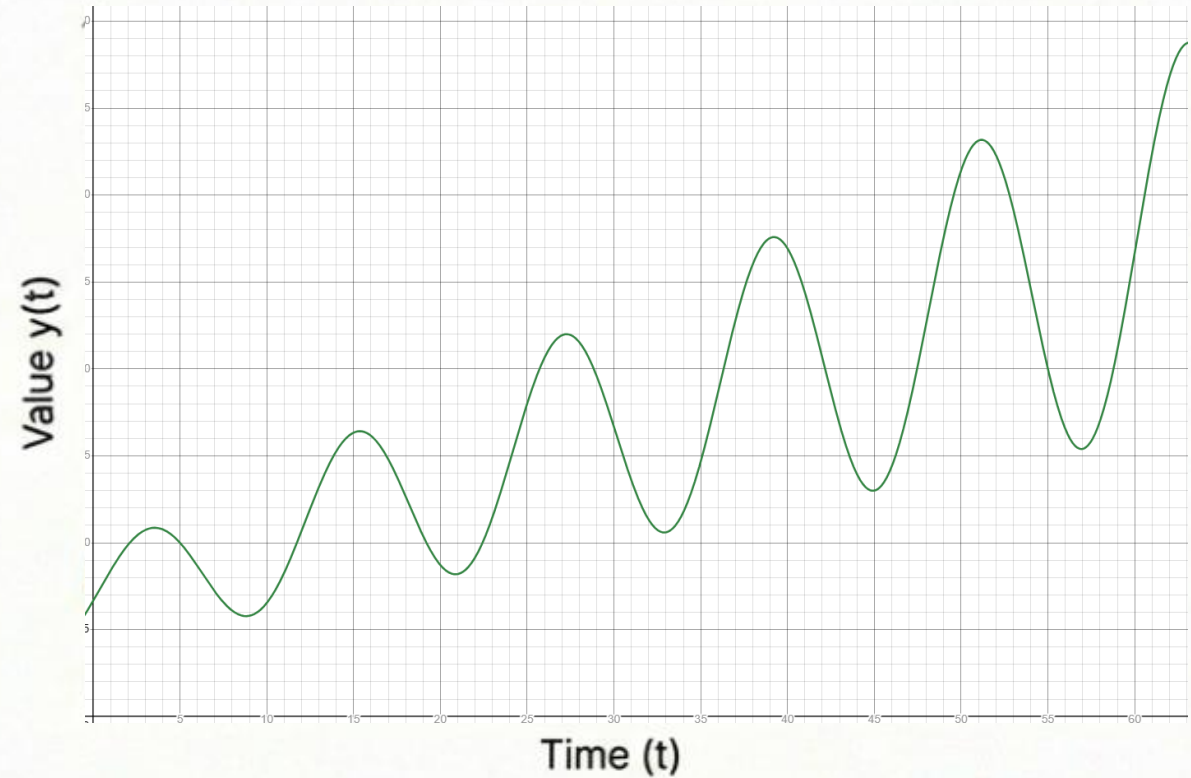
Window L

### Failure Mode

Window L

If Window < Period, the model is blind to the pattern.

**Additive**

$$y(t) = \text{Trend} + \text{Seasonality} + \text{Noise}$$

**Multiplicative**

$$y(t) = \text{Trend} \times \text{Seasonality} \times \text{Noise}$$

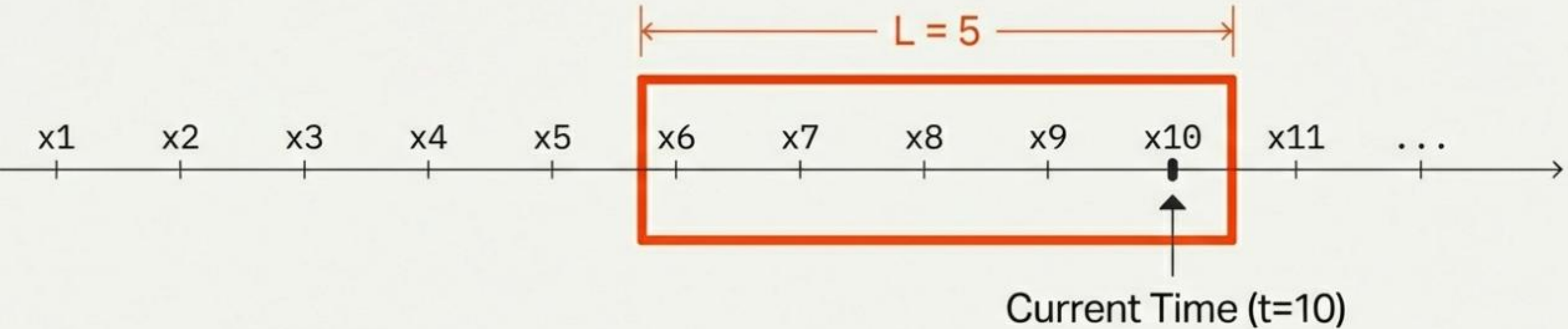E.g. Electricity usage, Temperature

E.g. Retail sales, Website traffic

# Windowing



Windowing is the process that converts the raw time series into a standard supervised learning dataset. By taking a specific chunk of the past, we create a context that asks the model to predict the next value.

# Lookback



The Lookback determines how far into the past the model sees at any given time step.

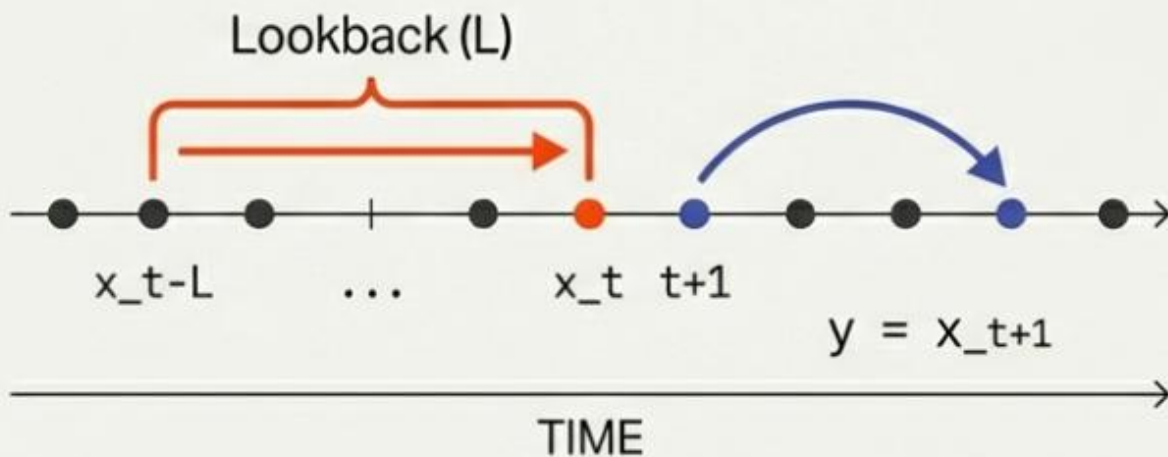Input Equation: Xt = [x_t-L+1, ..., x_t]

**EXAMPLE SCENARIO:**

Lookback L = 5
Time t = 10
Resulting Input Vector: [x6, x7, x8, x9, x10]

# Horizon

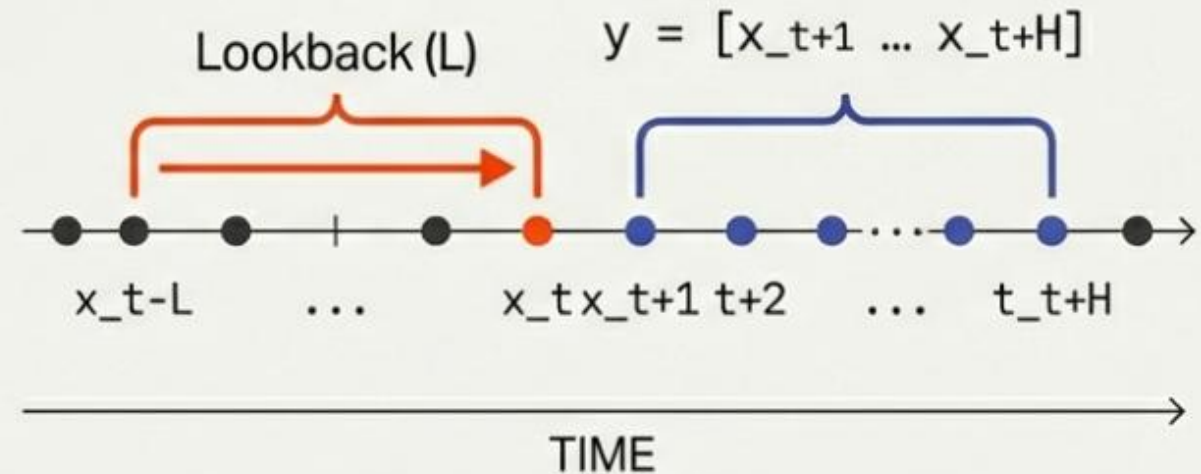The target (y) is what the model attempts to predict based on the Lookback.



## 1. One-step Target

Lookback (L)

x_t-L   ...   x_t   t+1

y = x_t+1

TIME

Use for immediate forecasts
(e.g., next hour).

## 2. Multi-step Target

Lookback (L)      y = [x_t+1 ... x_t+H]

x_t-L   ...   x_t x_t+1 t+2   ...   t_t+H

TIME

Use for forecast ranges
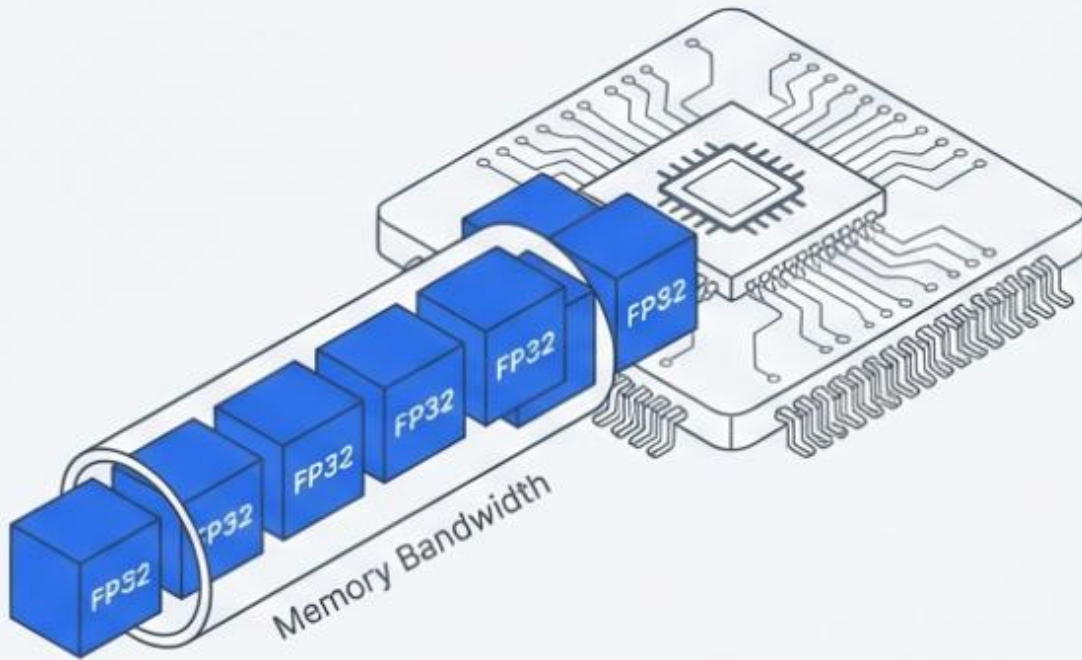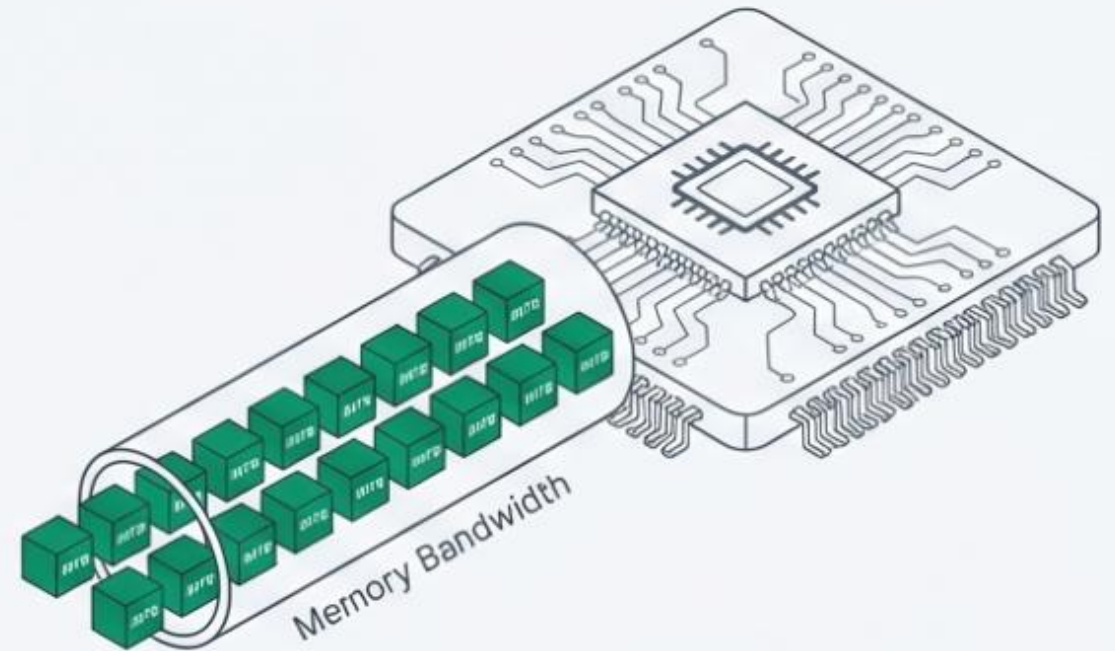(e.g., next 24 hours).

# Lab
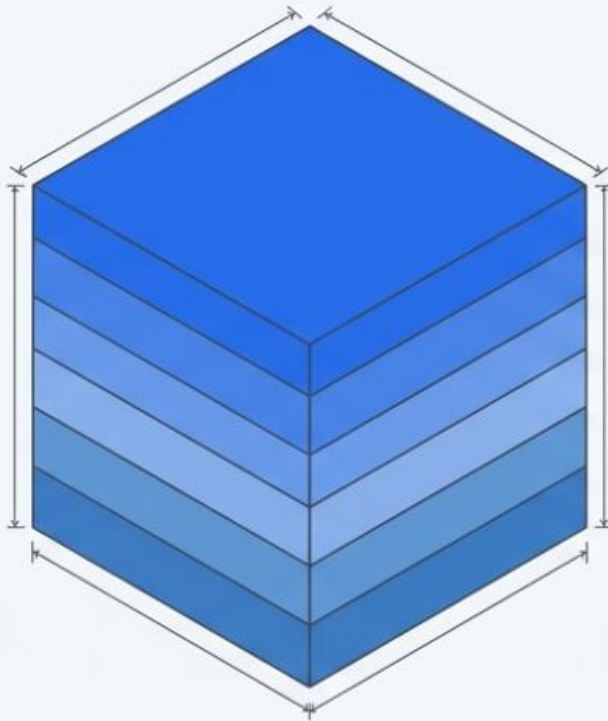
# Challenges in Production



## The Problem: Heavy Load

- **Compute Bound:** Massive multiply-add operations in Convolutions.
- **Memory Bound:** Moving 32-bit weights consumes high energy and time.
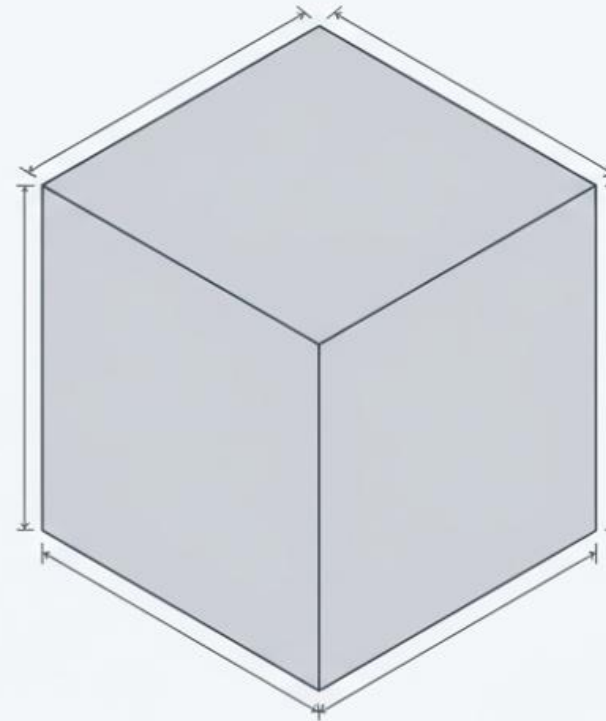
## The Solution: Quantization

- **Reduced Footprint:** 4x reduction in memory traffic (32-bit → 8-bit).
- **Accelerated Kernels:** Integer arithmetic executes significantly faster on edge CPUs.

# Weights (Per-Channel)



- **Applied to:** Model Parameters
- **Strategy:** Independent scale/zero-point per output channel.
- **Benefit:** High accuracy, accommodates varying ranges.

# Activations (Per-Tensor)



- **Applied to:** Inputs & Outputs
- **Strategy:** Single scale/zero-point for the entire layer.
- **Benefit:** Performance efficiency for dynamic data.

# Quantize and De quantize

## 01. QUANTIZE (Float → Int)

Discretizer

$$q = \text{clamp}(\text{round}(x / s) + z)$$

Range Limiter

## 02. DEQUANTIZE (Int → Float Approx)

$$\hat{x} \approx s * (q - z)$$

### Function Key

round(.) : Nearest integer conversion

clamp(.) : Restricts to [q_min, q_max]

q_min/max : Typically -128 to 127

**Range Inputs**

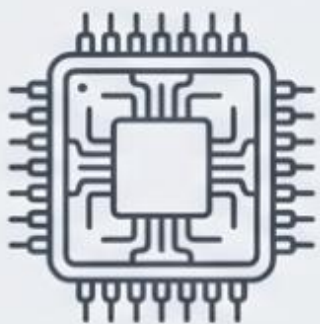$x\_min, x\_max$
$q\_min, q\_max$

**Scale Calculation (s)**

$$s = \frac{x\_max - x\_min}{q\_max - q\_min}$$

Spreads float range over integer codes.

**Zero-Point Calculation (z)**

$$z = \text{round}\left(q\_min - \frac{x\_min}{s}\right)$$

Aligns integer code z with float 0.

**Output**

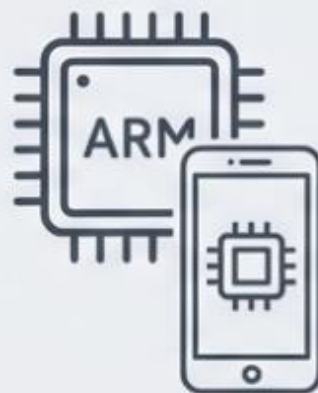Parameters s and z ready for quantization.

PyTorch 'torch.ao.quantization' relies on specific kernel backends.

## Server / Desktop

- Architecture: x86
- Backend Engine: `fbgemm` (Facebook GEMM)
- Optimized for: High throughput server-side inference.

## Mobile / Edge

- Architecture: ARM
- Backend Engine: `qnnpack` (Quantized Neural Network PACkage)
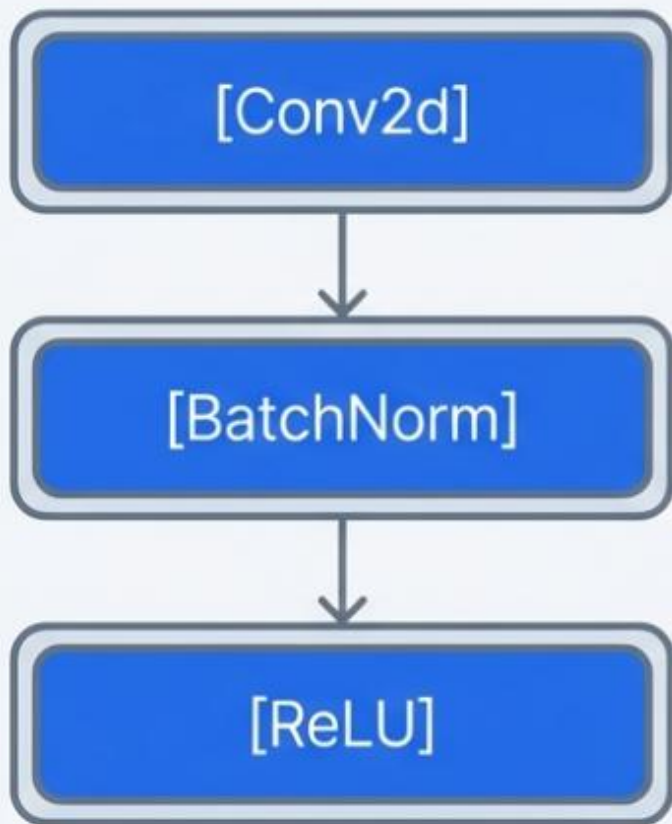- Optimized for: Low power, mobile processors.

Note: The theoretical concepts in this guide apply to both, but specific kernel support may vary.

# Fusing Operations

**Merging sequential operations for kernel efficiency.**

**Before Fusion**

[Conv2d]

↓

[BatchNorm]

↓

[ReLU]

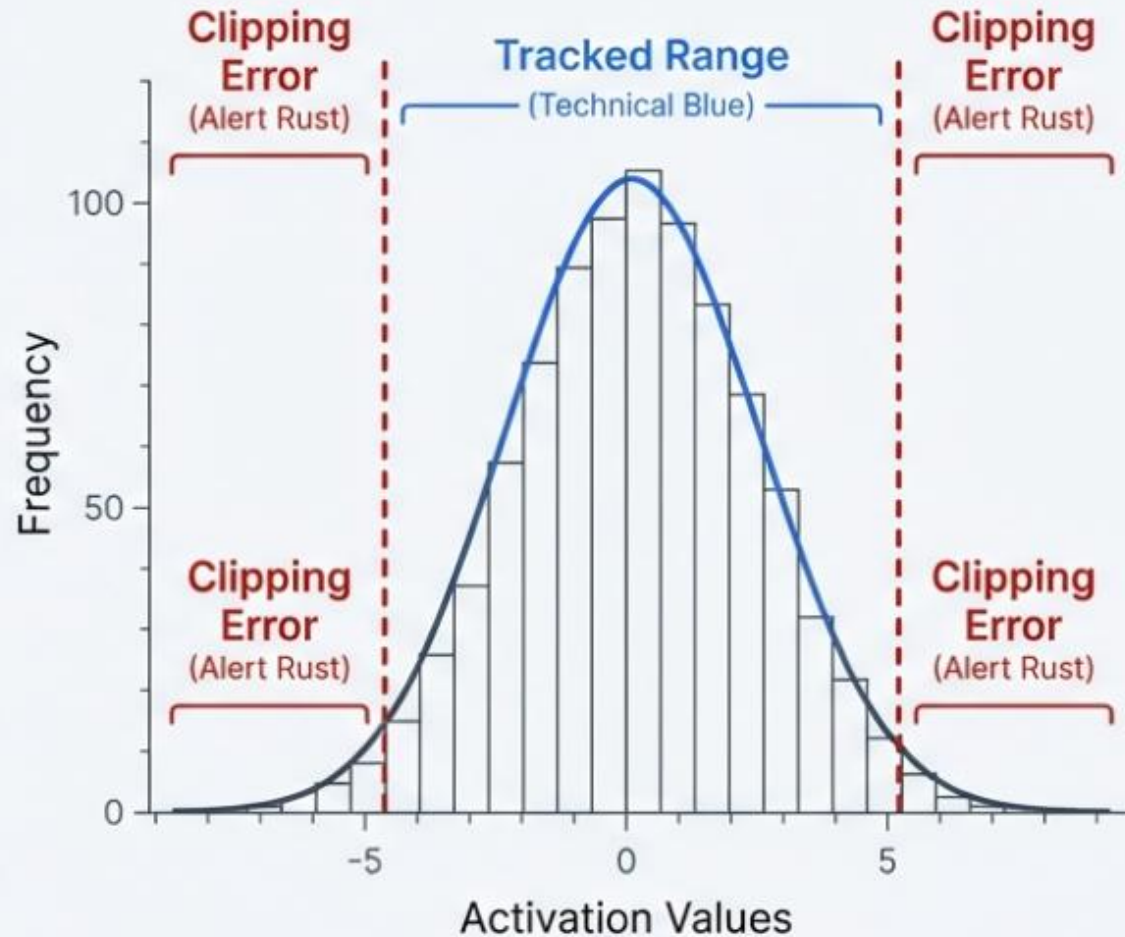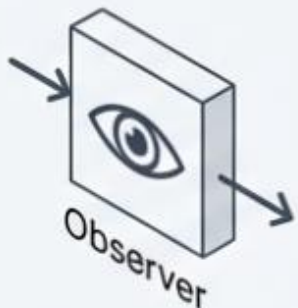**After Fusion**

[Conv2d + BN + ReLU]

**Why Fuse?**

- **Enables Quantization:** Backends like fbgemm require fused operators to access INT8 implementations.

- **Performance:** Reduces memory access overhead. Calculations happen in registers without writing intermediate results to VRAM.

# Quantizing Activations



**The Observer**

A module inserted into the network graph.

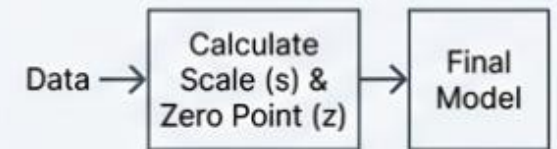It passively records statistics (Min/Max) of the tensors passing through it.

Observer

**The Calibration**
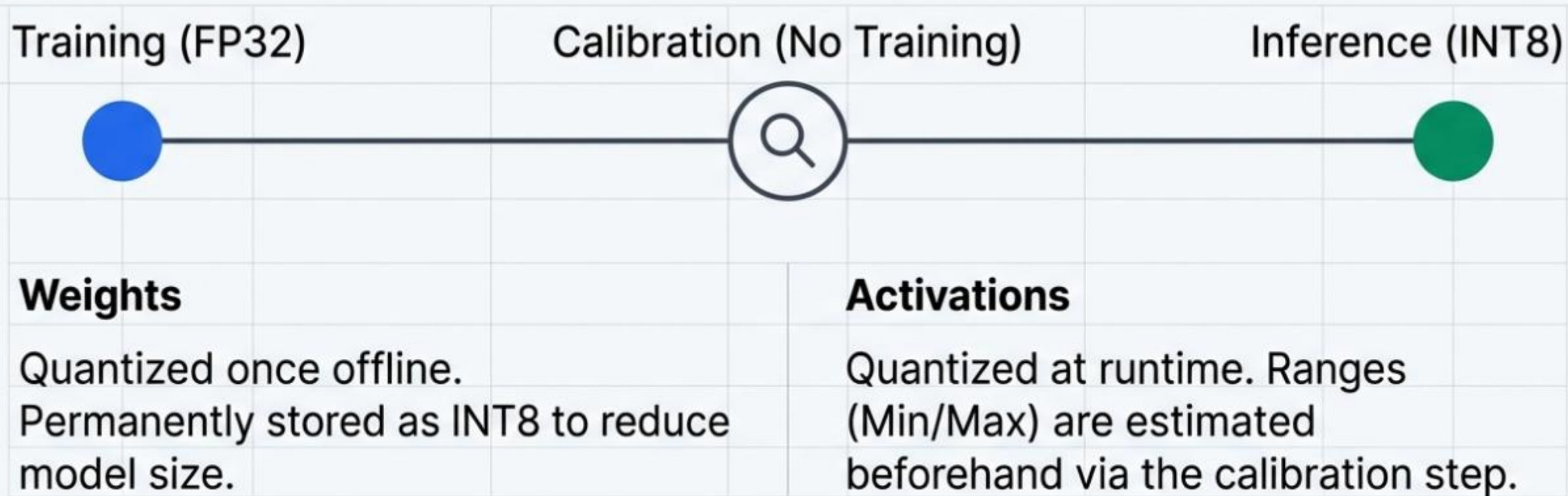
We run a forward pass with representative data.

No backpropagation occurs.

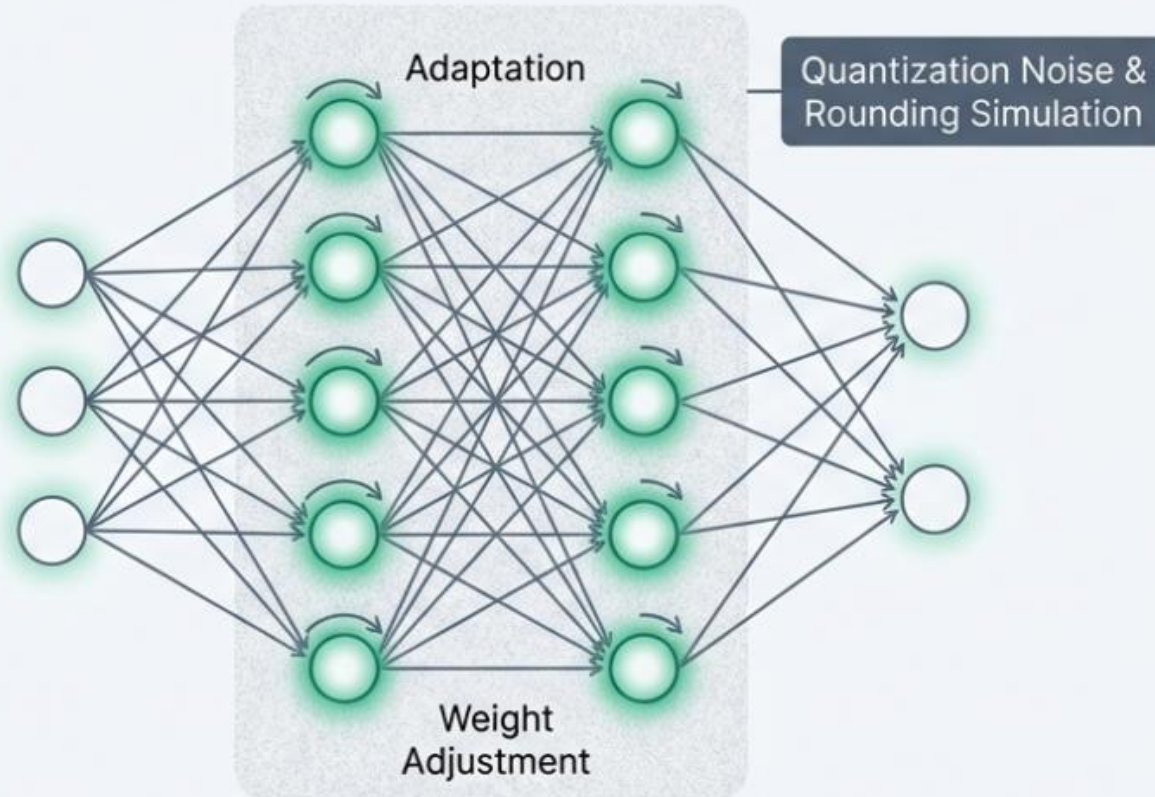The system uses these statistics to calculate the static Scale (s) and Zero Point (z) for the final model.

Data → Calculate Scale (s) & Zero Point (z) → Final Model

---

Clipping Error (Alert Rust)

Tracked Range (Technical Blue)

Clipping Error (Alert Rust)

Clipping Error (Alert Rust)

Clipping Error (Alert Rust)

Frequency

100

50

0

-5    0    5

Activation Values

A technique to quantize a pre-trained FP32 model without further backpropagation training steps.

| Training (FP32) | Calibration (No Training) | Inference (INT8) |
|---|---|---|

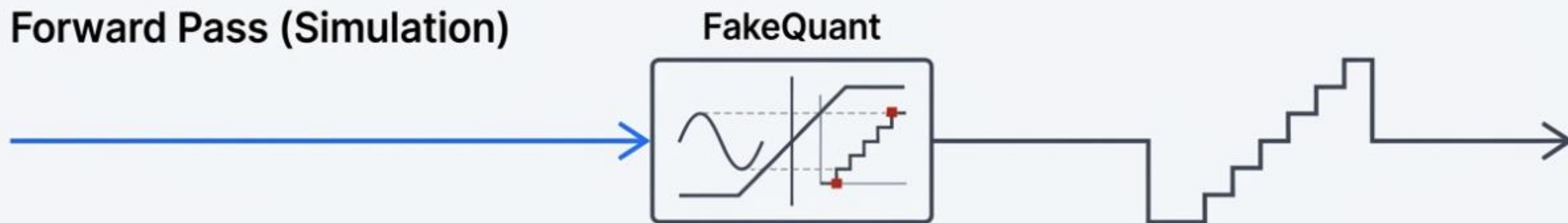| Weights | | Activations |
|---|---|---|
| Quantized once offline. Permanently stored as INT8 to reduce model size. | | Quantized at runtime. Ranges (Min/Max) are estimated beforehand via the calibration step. |

# Quantization Aware Training



**Core Concept: Adaptation.**

- Instead of quantizing AFTER training (blindly), we simulate quantization DURING training.

- The network learns to adjust its weights to survive the noise and rounding errors of INT8 representation.

# Forward Pass (Simulation)

## FakeQuant



Simulates rounding/clamping. The loss function "sees" the quantization error.

# Backward Pass (Learning)

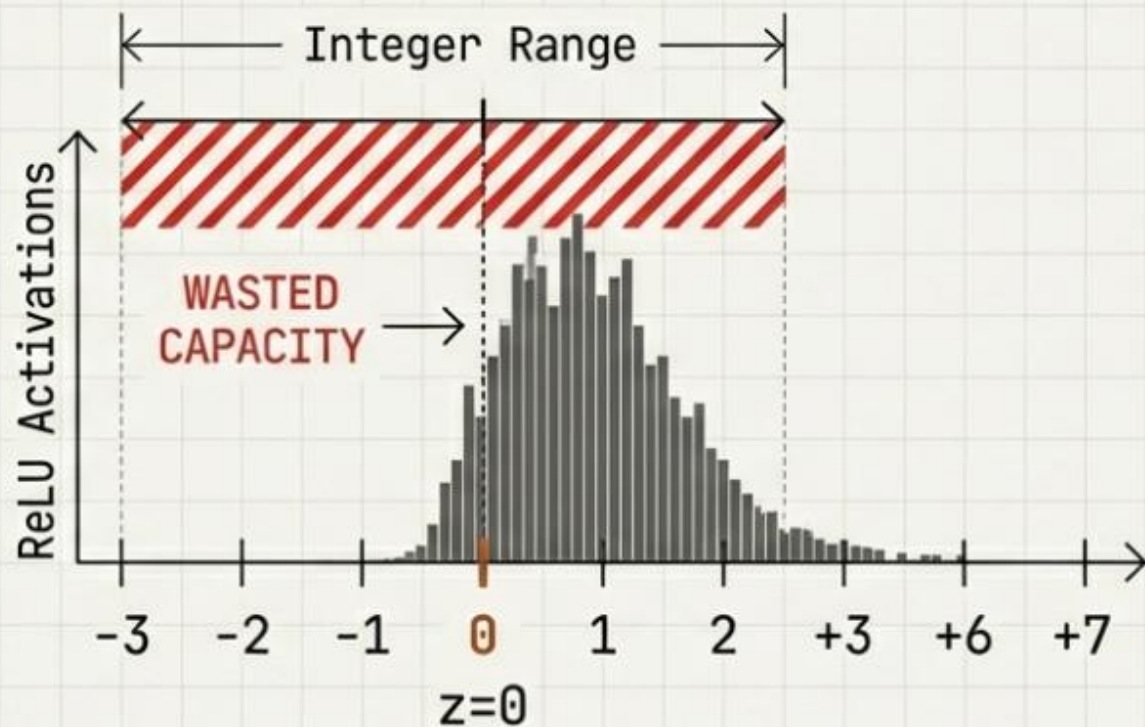Gradients update the high-precision FP32 weights. The master weights cluster into robust values.
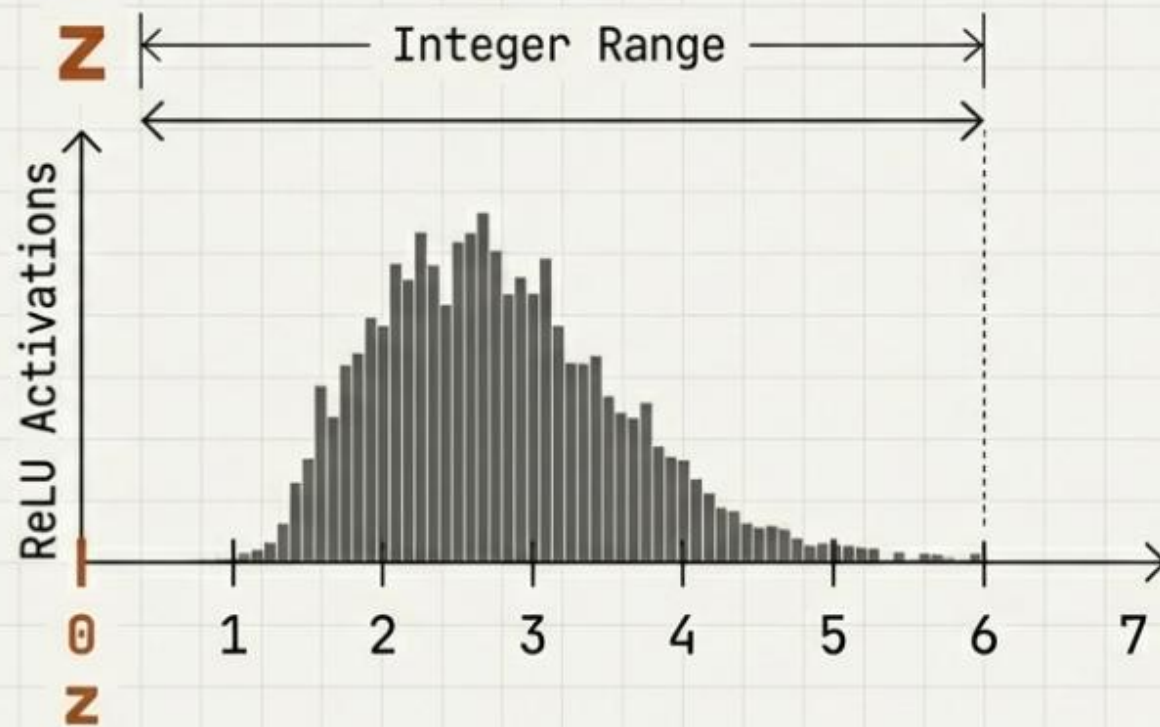
# Lab

# Thank You

# Appendix