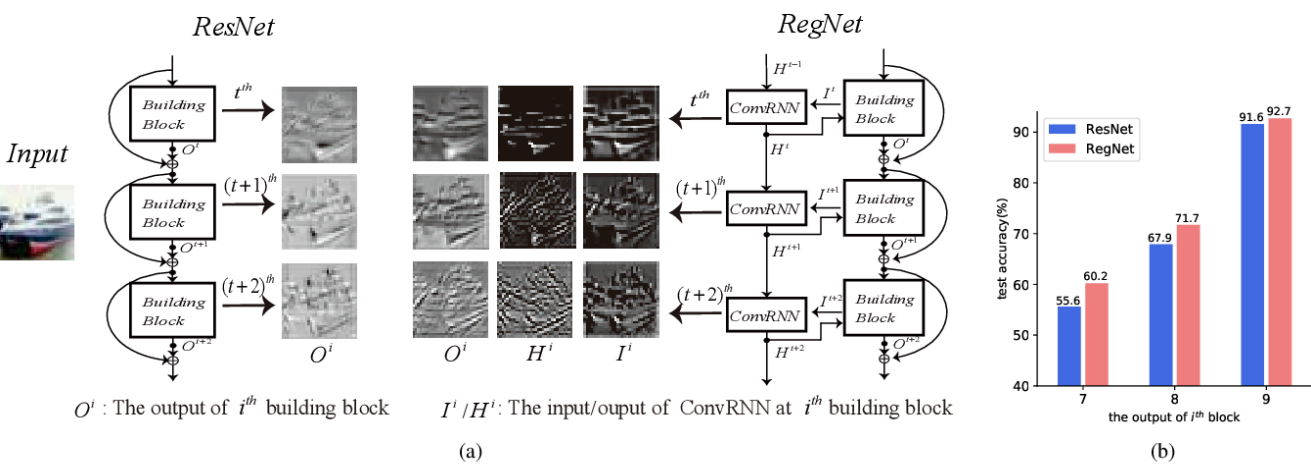# Research Paper Notes

The paper introduces a paradigm shift from **Neural Architecture Search (NAS)**—which optimizes for a single best model instance—to **Design Space Design**, which optimizes a parametrization for a whole population of models. This approach aims to discover general design principles that are robust and interpretable, rather than just a specific architecture tuned for one setting.



$O^i$ : The output of $i^{th}$ building block          $I^i/H^i$: The input/ouput of ConvRNN at $i^{th}$ building block

(a)                                                                          (b)

---

# 1. Comparison to Existing Methods in the Context of RegNet

## Manual Network Design

| Aspect | Manual Network Design | Designing Network Design Spaces (RegNet) |
|---|---|---|
| **Focus/Outcome** | Discovery of new design choices and generalized design principles (e.g., LeNet, ResNet). | Discovery of **general design principles** but achieved at the **design space level**. |
| **Methodology** | Largely manual process focused on improving accuracy. | Analogous to manual design but **elevated to the population level** and guided by distribution estimates. |
| **Limitation Addressed** | Finding well-optimized networks manually becomes challenging as the number of design choices increases. | Uses semi-automated procedures and rigorous population analysis (EDF) to refine and simplify the design space. |

## Neural Architecture Search (NAS)

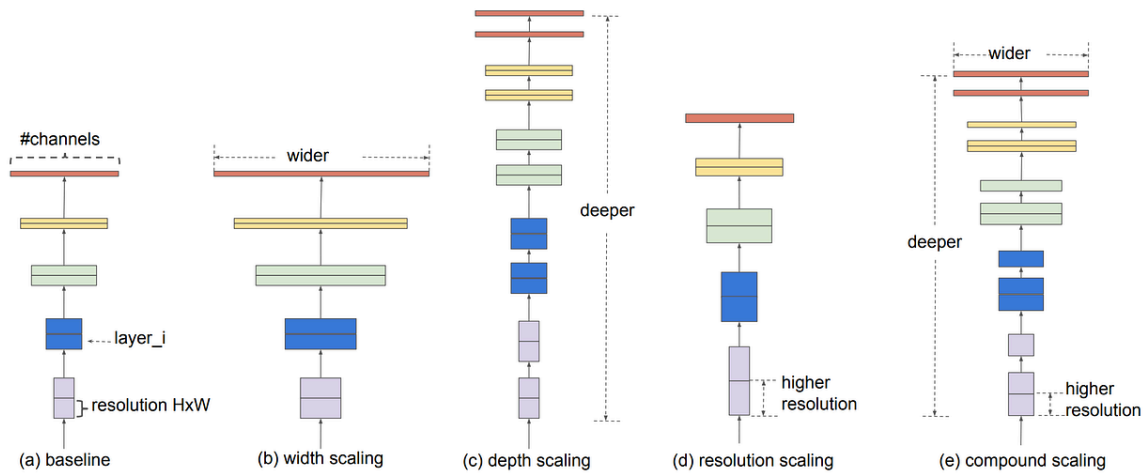| Aspect | Neural Architecture Search (NAS) | Designing Network Design Spaces (RegNet) |
|---|---|---|
| **Focus/Outcome** | Efficiently finding the **best network instance** tuned to a specific setting (e.g., hardware platform). | Designing the **design space itself** to parametrize populations of networks and discover general principles. |
| **Limitations** | Does not enable the discovery of network design principles that deepen understanding and allow generalization to new settings. | Aims to find **simple models that are easy to understand, build upon, and generalize**. |
| **Relationship** | The two approaches are **complementary**: better design spaces, like RegNet, can improve the efficiency of NAS search algorithms and lead to the existence of better models by enriching the design space. | |

# 2. Methodology: AnyNet to RegNet Refinement

## Summary

The Design Space Design methodology progressively simplifies a massive, unconstrained search space (AnyNet) into a structured one (RegNet) by analyzing model populations. By applying constraints like shared parameters and increasing dimensions, researchers reduced $10^{18}$ configurations to a simple linear model without losing accuracy.

## Key Points

- **Goal:** Simplify the search space while improving interpretability and maintaining model diversity.
- **Primary Tool (EDF):** Uses the **Error Empirical Distribution Function** to analyze the quality of the *entire population*, not just the best model.
- **Proxy Training:** Uses low-compute training (e.g., 10 epochs) to efficiently evaluate millions of sampled models.
- **Reduction:** Successfully reduced the design space size by 10 orders of magnitude (from $\sim 10^{18}$ to $\sim 10^{8}$).

(a) baseline  (b) width scaling  (c) depth scaling  (d) resolution scaling  (e) compound scaling

# Step-by-Step Explanation

The refinement process moved from a chaotic space to a highly ordered one through 5 specific steps:

1. **AnyNetXA (Base):** Unconstrained. Standard ResNet-like blocks where width, depth, and groups vary freely. (16 degrees of freedom).
2. **AnyNetXB (Shared Bottleneck):** Constraint applied: $b_i = b$. All stages share the same bottleneck ratio. Quality maintained.
3. **AnyNetXC (Shared Groups):** Constraint applied: $g_i = g$. All stages share the same group width. Space shrinks by $10^4$.
4. **AnyNetXD (Increasing Widths):** Pattern observed: Good networks get wider. Constraint applied: $w_{i+1} \geq w_i$. **Error improves.**
5. **AnyNetXE (Increasing Depths):** Pattern observed: Good networks get deeper. Constraint applied: $d_{i+1} \geq d_i$. **Error improves.**
6. **RegNet (Linear Fit):** Analysis showed optimal width growth is linear. Replaced free variables with a quantized linear function: $u_j = w_0 + w_a \cdot j$.

---

# Example

**The Mathematical Simplification**

Instead of manually selecting the width for every single block (which is hard to optimize), RegNet proves you only need a linear slope equation.

- AnyNet (Unconstrained): Layer_Widths = [48, 120, 64, 256, 128...] (Hard to search)
- RegNet (Linear Constraint): Layer_Widths = [64, 128, 192, 256...] (Defined by Slope $w_a$)

This linear structure is defined by:

$$u_j = w_0 + w_a \cdot j$$

(Where $w_0$ is starting width and $w_a$ is the slope)

# 3. RegNet: Design Principles & Key Findings

## Summary

RegNet simplifies neural network design by proving that the optimal width and depth of layers follow a **quantized linear function**, reducing the design space from 16 complex parameters to just 6 simple ones. Analysis of RegNet reveals surprising principles: optimal networks often have stable depth (~20 blocks) and no bottlenecks ($b = 1.0$), making them significantly faster on GPUs than EfficientNet.

## Key Points

- **Quantized Linear Function:** The core insight is that network width should grow linearly with depth ($u_j = w_0 + w_a \cdot j$).
- **Dimensionality Reduction:** Simplified the design space from 16 parameters (AnyNetX) to 6 (RegNet), shrinking the space by 10 orders of magnitude.
- **Stable Depth:** Unlike common practice, the best models stay at a stable depth of ~20 blocks (60 layers) even as compute increases.
- **No Bottlenecks:** The best performing models utilize a bottleneck ratio of $b = 1.0$, effectively removing the bottleneck.
- **GPU Speed:** RegNet is up to **5× faster** than EfficientNet because its activations scale with the square root of FLOPs ($\sqrt{\text{flops}}$) rather than linearly.

## Example

**RegNet Wisdom vs. Common Practice**

| Feature | Common Practice (e.g., EfficientNet/MobileNet) | RegNet Finding |
|---|---|---|
| **Depth** | Deeper is better for high compute. | **Stable depth** (~20 blocks) is best across regimes. |
| **Bottlenecks** | Use Inverted Bottlenecks ($b < 1$) for efficiency. | Remove bottlenecks ($b = 1.0$) for performance. |
| **Resolution** | Scale input resolution (e.g., 600px+). | **Fixed resolution** (224x224) is sufficient. |
| **Scaling** | Double width at every stage ($w_m = 2$). | Use a multiplier of $w_m \approx 2.5$. |

# 4. RegNet: Performance Benchmarks & Comparisons

## Summary

RegNet demonstrates superior performance over state-of-the-art models like EfficientNet and ResNe(X)t by optimizing the design space rather than individual instances. Notably, RegNet is up to **5× faster on GPUs** due to efficient activation scaling ($\sqrt{\text{flops}}$), while maintaining high accuracy across mobile and high-compute regimes.

## Key Points

- **GPU Speed:** RegNet is up to 5× faster than EfficientNet because its activations scale with the square root of FLOPs, not linearly.
- **Fair Comparison:** Outperforms EfficientNet and ResNe(X)t when trained under identical, controlled conditions (no extra regularization).
- **Mobile Efficiency:** RegNetY-600MF outperforms MobileNet-V2 and NASNet-A using only a basic training schedule.
- **Resolution:** Unlike EfficientNet, RegNet performs best with a **fixed resolution** ($224 \times 224$), even at higher FLOPs.
- **Variants: RegNetY** (RegNetX + Squeeze-and-Excitation) consistently improves performance.

**Understanding the Performance Gap**

1. **Controlled Training:** Researchers stripped away "tricks" (like AutoAugment) to compare architectures purely on design quality.
2. **Activation Scaling:**
   - **EfficientNet:** Activations scale linearly with FLOPs (due to increasing resolution/depth). This bottlenecks GPUs.
   - **RegNet:** Activations scale with $\sqrt{\text{flops}}$. This allows much faster training and inference on accelerators.
3. **Regime Analysis:**
   - *Low Compute:* EfficientNet is competitive.
   - *High Compute:* RegNet pulls ahead in both accuracy and speed.
4. **Design Choice Validation:** Tests confirmed that **inverted bottlenecks** ($b < 1$) hurt performance, while **standard bottlenecks** ($b = 1$) are optimal.
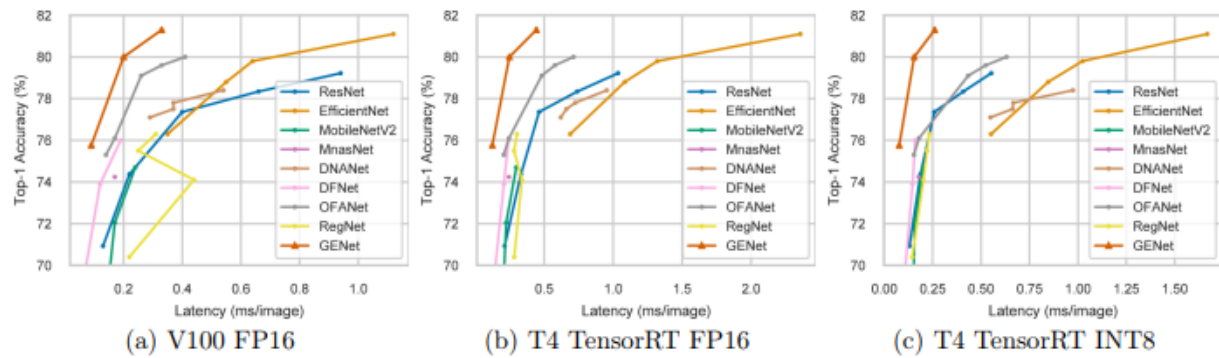
Figure 5: ImageNet top-1 accuracy v.s. network latency, batch size 64 on V100, batch size 32 on T4.

# Example

**The "RegNetX-8.0GF" Speed Scenario**

When comparing high-compute models under identical conditions:

- **Model A:** EfficientNet-B5
- **Model B:** RegNetX-8.0GF

Result:

RegNetX-8.0GF achieves lower error and runs 5× faster than EfficientNet-B5. This makes it viable for real-time applications where EfficientNet would be too slow.