
Contextual Feedback Loops: Amplifying Deep Reasoning with Iterative Top-Down Feedback

Jacob Fein-Ashley¹

Abstract

Deep neural networks typically rely on a single forward pass for inference, which can limit their capacity to resolve ambiguous inputs. We introduce **Contextual Feedback Loops** (CFLs) as an iterative mechanism that incorporates top-down feedback to refine intermediate representations, thereby improving accuracy and robustness. This repeated process mirrors how humans continuously re-interpret sensory information in daily life—by checking and re-checking our perceptions using contextual cues. Our results suggest that CFLs can offer a straightforward yet powerful way to incorporate such contextual reasoning in modern deep learning architectures.

1. Introduction

Deep learning architectures, ranging from convolutional neural networks to transformers, have driven remarkable advances in image recognition, natural language processing, and other domains (LeCun et al., 2015). Despite their impressive capabilities, many of these successes rely on feed-forward processing: a single pass of information from lower-level features to higher-level abstractions. While effective for numerous tasks, such strictly bottom-up pipelines can struggle with intricate or ambiguous inputs that naturally benefit from iterative reflection and contextual reinterpretation.

Human perception, by contrast, often relies on top-down guidance. In everyday life, people continuously refine their senses by leveraging higher-level expectations. For instance, when trying to identify a distant figure on a foggy street, we might guess it is a neighbor and then *re-check* the silhouette, adjusting our perception if it doesn't match that top-down guess. This everyday reasoning loop, where expectations reshape early perception, stands in stark contrast to the single-pass nature of conventional deep neural networks.

We propose **Contextual Feedback Loops** (CFLs) as a mechanism to incorporate this form of top-down feedback within modern neural architectures. Instead of halting at the initial forward pass, CFLs iteratively revisit intermediate represen-

tations, guided by a high-level context vector derived from the model's own predictions. Through these repeated inference steps, the network can potentially refine ambiguous features in light of emerging semantic cues, bridging the gap between purely bottom-up computation and the iterative interpretation observed in human cognition.

To illustrate this process, consider a speech recognition scenario where background noise or overlapping voices initially obscure the audio. By forming a hypothesis of which words might be spoken, the model can re-weight certain frequency bands or refine its acoustic representation to better match the context. CFLs automate this form of hypothesis-driven perception, offering a path toward more context-aware and interpretable learning processes in real-world settings.

The remainder of this paper is structured as follows:

- **Section 2 (Related Work)** reviews core ideas from cognitive science, predictive coding, and existing feedback-based neural frameworks, situating CFLs in a broader historical and theoretical context.
- **Section 3 (Methods)** formalizes the concept of Contextual Feedback Loops, detailing the mathematical framework, the iterative refinement procedure, and the backpropagation-through-time training approach.
- **Section 4 (Experiments)** evaluates CFLs on benchmark tasks, comparing their performance to feed-forward baselines and demonstrating tangible gains in classification accuracy and convergence.
- **Section 5 (Conclusion)** summarizes our contributions and outlines future directions, including adapting CFLs to more diverse architectures and larger-scale datasets.

Taken together, these contributions aim to show how even modest feedback loops can enhance a neural network's ability to interpret and refine its internal representations, moving closer to the dynamic, context-rich reasoning that shapes human perception.

2. Related Work

Integrating top-down feedback into neural architectures is deeply rooted in cognitive science and neuroscience. Cognitive theories such as Adaptive Resonance Theory (ART) (Grossberg, 2017) highlight the importance of resonant feedback loops in stabilizing and refining perceptual states. ART suggests that conscious perception emerges from a match between top-down predictions and bottom-up sensory signals, resonating with the motivation for Contextual Feedback Loops (CFLs), where high-level representations guide the re-interpretation of earlier-layer activations.

From a computational perspective, predictive coding frameworks posit that perception arises from the interaction between top-down generative models and bottom-up error signals (Rao & Ballard, 1999; Friston, 2010). These theories have inspired numerous models that blend inference and learning, including Bayesian approaches to visual perception (Lee & Mumford, 2003; Yuille & Kersten, 2006) and deep generative models that leverage iterative refinements to improve latent representations (Hinton et al., 1995; Rezende et al., 2014; Greff et al., 2019).

Recurrent and feedback mechanisms in artificial neural networks have shown improvements in object recognition, scene understanding, and video prediction tasks by iteratively refining internal states (Spoerer et al., 2017; Lotter et al., 2017; Wen et al., 2018). Similarly, work on bidirectional learning (Adigun & Kosko, 2020) and inverse mappings in neural networks shows that integrating backward passes for reasoning and interpretation can enhance robustness and offer insights into network behavior.

Other studies have also explored how incorporating top-down or synthetic gradient signals can accelerate learning and stabilization (Jaderberg et al., 2017), while hierarchical models with iterative inference steps, such as those based on variational autoencoders (Kingma & Welling, 2014) and iterative refinement techniques (Andrychowicz et al., 2016), have demonstrated that looping feedback through a network can yield more compact and disentangled representations.

Brain-Inspired Predictive Coding Dynamics. Recent work has explicitly explored predictive coding within deep neural networks. For instance, Choksi et al. (Choksi et al., 2021) proposed *Predify*, which augments feed-forward DCNNs with generative feedback connections trained via an unsupervised reconstruction objective. Across multiple timesteps, these feedback signals reduce reconstruction error and improve robustness to perturbations. While *Predify* emphasizes layer-specific generative feedback, our CFL framework more generally re-injects high-level context for iterative refinement, supporting diverse architectures with flexible adapters.

Comparison with Feedback Networks. Zamir et al. (Zamir et al., 2017) introduced *Feedback Networks*, which unroll a recurrent module over multiple timesteps in a feed-forward backbone. Although both Feedback Networks and CFLs use iterative updates, CFL explicitly feeds output-derived context back to earlier layers, enabling top-down cues to guide refinements more directly. In contrast, Feedback Networks primarily deepen representations via recurrent transformations. CFL thus functions as a general mechanism for context-driven reinterpretation, integrating seamlessly with CNNs or Transformers while preserving interpretability and efficient training.

In summary, CFLs build on a tradition emphasizing top-down and bottom-up interplay. By offering a principled strategy to repeatedly inject output-level context into earlier representations, CFLs align with biologically inspired, generative, and recurrent modeling approaches, resulting in more context-aware, interpretable, and robust deep learning systems.

3. Methods

3.1. Motivation

Many neural network architectures, including feed-forward, convolutional, recurrent, or transformer-based models, process information primarily in a single forward direction, generating representations that flow from lower-level features to higher-level abstractions. While effective in many scenarios, this paradigm can be constrained when the input or task is inherently ambiguous or requires iterative reasoning and contextual interpretation.

Drawing inspiration from human perception, where top-down signals can influence and refine lower-level feature processing, we introduce **Contextual Feedback Loops** (CFLs). These loops provide a mechanism to integrate high-level contextual information back into earlier stages of the network. The goal is to enable a wide range of neural architectures to iteratively refine their internal representations, thereby improving robustness and interpretability when faced with complex or ambiguous inputs.

3.2. Usefulness in Context-Rich Tasks

Context often provides additional semantic cues that can substantially influence the interpretation of raw data. For instance, a subtle shadow might distinguish a familiar face from its surroundings, or a slight intonation may reveal the sentiment behind a spoken phrase. In many real-world problems, the relationship between high-level semantics and low-level observations is complex and cannot be adequately captured by a purely feed-forward mechanism.

Contextual Feedback Loops address this limitation by en-

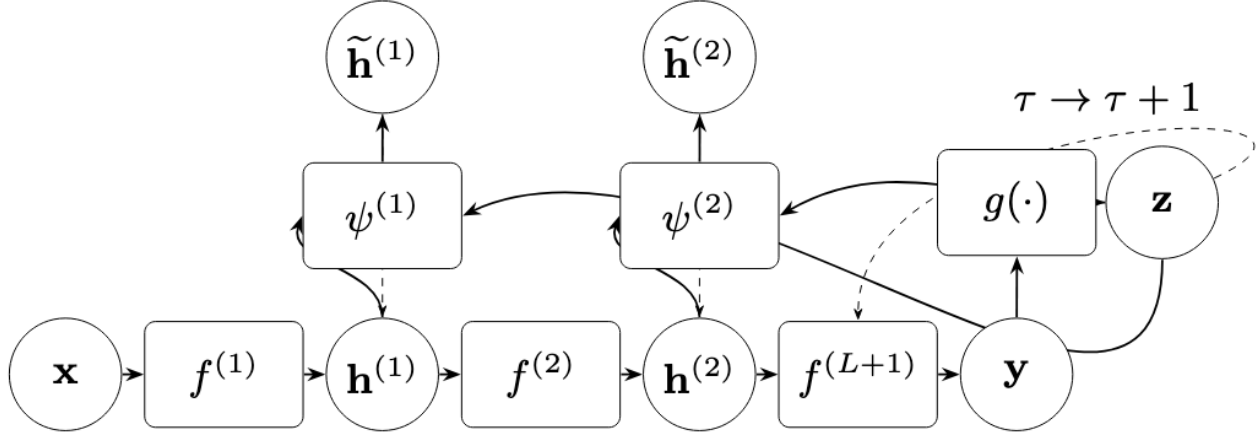


Figure 1. An overview of the Contextual Feedback Loops (CFL) framework. The top-level context vector is derived from the output and fed back to earlier layers, enabling iterative refinement.

abling a model to iteratively refine its hidden representations in light of the best current estimate of the output. Rather than treating top-level predictions as a static final result, CFL continuously re-injects these predictions back into earlier layers. This iterative feedback mechanism helps reconcile potential inconsistencies between top-down expectations and bottom-up observations, leading to more coherent and context-aware representations.

3.3. General Framework

Consider a generic neural network that, given an input $\mathbf{x} \in \mathbb{R}^{d_x}$, produces an output $\mathbf{y} \in \mathbb{R}^{d_y}$. The network can be composed of multiple layers or modules arranged in any architecture (e.g., feed-forward stack, convolutional layers, attention blocks, recurrent cells), forming an overall differentiable function:

$$\mathbf{y} = \mathcal{F}(\mathbf{x}; \theta), \quad (1)$$

where θ represents all learnable parameters of the network.

In a standard setting, the network performs a single forward pass from \mathbf{x} to \mathbf{y} . To incorporate iterative refinement, we introduce an additional pathway for top-down context to influence the intermediate representations. Let $\{\mathbf{h}^{(l)}\}$ denote these intermediate states for layers or components $l \in \{1, \dots, L\}$, so that:

$$\mathbf{h}^{(1)} = f^{(1)}(\mathbf{x}), \quad \mathbf{h}^{(2)} = f^{(2)}(\mathbf{h}^{(1)}), \dots, \quad \mathbf{y} = f^{(L+1)}(\mathbf{h}^{(L)}), \quad (2)$$

with each $f^{(l)}(\cdot)$ representing a portion of the network.

3.4. Feedback Integration

To enable top-down feedback, we define a mapping from the output (or a high-level representation close to the output)

back into a **context vector**:

$$\mathbf{z} = g(\mathbf{y}), \quad (3)$$

where $g(\cdot)$ is a learned transformation. This context vector \mathbf{z} encapsulates high-level semantic information derived from the network’s output. It is integrated back into intermediate layers to influence their representations.

To incorporate \mathbf{z} into intermediate layers, we define a set of **feedback adapters** $\{\psi^{(l)}\}$, each of which takes the current hidden state $\mathbf{h}^{(l)}$ and the context \mathbf{z} as input and produces a refined representation:

$$\tilde{\mathbf{h}}^{(l)} = \psi^{(l)}(\mathbf{h}^{(l)}, \mathbf{z}). \quad (4)$$

These adapters can be implemented through linear gating functions, attention mechanisms, or other learnable transformations. Their form is flexible and can be adapted to various network architectures.

3.5. Iterative Refinement Procedure

The core idea of CFL is to alternate between forward computation of the output and top-down refinement of intermediate representations. Let $\tau = 0, 1, \dots, T$ index the refinement steps.

1. **Initialization (Forward Pass):** At $\tau = 0$, perform a forward pass through the network to obtain the initial output:

$$\mathbf{y}^{(0)} = \mathcal{F}(\mathbf{x}; \theta). \quad (5)$$

2. **Context Computation:** Compute the top-down context vector:

$$\mathbf{z}^{(\tau)} = g(\mathbf{y}^{(\tau)}). \quad (6)$$

3. **Refinement of Hidden States:** Each intermediate representation is updated to incorporate the context:

$$\mathbf{h}_{\tau+1}^{(l)} = \alpha \mathbf{h}_{\tau}^{(l)} + (1 - \alpha) \psi^{(l)}(\mathbf{h}_{\tau}^{(l)}, \mathbf{z}^{(\tau)}), \quad (7)$$

where $\alpha \in [0, 1]$ controls how strongly the updated state replaces the original.

4. **Recompute Output:** After updating all $\mathbf{h}_{\tau}^{(l)}$, pass the refined representations forward:

$$\mathbf{y}^{(\tau+1)} = f^{(L+1)}(\mathbf{h}_{\tau+1}^{(L)}). \quad (8)$$

5. **Repeat Until Convergence or Max Steps:** Repeat the context computation and refinement steps until T iterations or until a convergence criterion is met.

The final refined output is taken as $\mathbf{y}^{(T)}$. This iterative process effectively integrates top-down context into the network’s intermediate states multiple times, allowing it to resolve ambiguities and refine its predictions.

3.6. Training via Backpropagation Through Time

The iterative refinement introduces a temporal dimension to the inference process. To train the parameters θ of the network, as well as those of $g(\cdot)$ and the adapters $\{\psi^{(l)}\}$, we can unroll the computation for T steps and apply backpropagation through time (BPTT).

Given training data $(\mathbf{x}, \mathbf{y}^*)$, where \mathbf{y}^* is the target output, we define a loss function over the sequence of predictions:

$$\mathcal{L} = \sum_{\tau=0}^T \lambda_{\tau} \ell(\mathbf{y}^{(\tau)}, \mathbf{y}^*), \quad (9)$$

where $\ell(\cdot)$ is a standard loss function (e.g., cross-entropy) and λ_{τ} weights the contribution of each refinement step.

Since all operations are differentiable, gradients flow through the iterative loops, enabling end-to-end training. Standard optimization methods (e.g., SGD, Adam) can then be used to update all trainable parameters.

3.7. Applicability to Different Architectures

This framework is not restricted to any particular class of neural networks. Any layered or modular architecture that provides access to intermediate states can incorporate CFL:

- **Convolutional networks:** Integrate \mathbf{z} into feature maps via gating or attention-based adapters.
- **Recurrent networks:** Incorporate context into hidden states at each refinement step to re-evaluate sequential information.

- **Transformer models:** Inject \mathbf{z} as an additional conditioning vector into attention blocks or feed-forward layers.

In each case, a top-down context vector is derived from the network’s output and reintroduced into intermediate states, guiding iterative refinement and promoting improved coherence in the final predictions.

3.8. Theoretical Discussion, Convergence Insights, and Proofs

From a theoretical standpoint, **Contextual Feedback Loops** can be interpreted as iteratively seeking a *fixed point* in the space of hidden representations. Formally, each layer’s hidden state update at step τ can be written as:

$$\mathbf{h}_{\tau+1}^{(l)} = F^{(l)}(\mathbf{h}_{\tau}^{(l)}, \mathbf{z}_{\tau}), \quad (10)$$

for some transformation $F^{(l)}$ that depends on both the feed-forward pathway and the top-down feedback adapter $\psi^{(l)}$. Over multiple refinement steps, the system attempts to converge to a point where

$$\mathbf{h}_{\tau+1}^{(l)} \approx \mathbf{h}_{\tau}^{(l)}, \quad (11)$$

simultaneously for all layers l .

Under mild assumptions (such as contractive mappings in the hidden-state space or a damping factor α that limits the magnitude of updates), one can analyze the stability of this iterative process. If each step is sufficiently small (i.e., α is chosen appropriately) and $\psi^{(l)}$ exhibits Lipschitz continuity, then fixed-point convergence can often be assured.

Contraction Mapping Perspective. Let $\mathbf{h} = [\mathbf{h}^{(1)}, \dots, \mathbf{h}^{(L)}]$ denote the concatenation of hidden states across all layers. Define a global transformation

$$F(\mathbf{h}) = [F^{(1)}(\mathbf{h}^{(1)}, \mathbf{z}), \dots, F^{(L)}(\mathbf{h}^{(L)}, \mathbf{z})], \quad (12)$$

where $\mathbf{z} = g(f^{(L+1)}(\mathbf{h}^{(L)}))$ is the top-down context derived from the current output.

Theorem 3.1 (Convergence under Contraction). *Suppose there exists a constant $0 \leq \gamma < 1$ such that for all pairs \mathbf{h}, \mathbf{h}' ,*

$$\|F(\mathbf{h}) - F(\mathbf{h}')\| \leq \gamma \|\mathbf{h} - \mathbf{h}'\|.$$

Then, by the Banach Fixed-Point Theorem, there exists a unique fixed point \mathbf{h}^ satisfying $F(\mathbf{h}^*) = \mathbf{h}^*$, and the iteration*

$$\mathbf{h}_{\tau+1} = \alpha \mathbf{h}_{\tau} + (1 - \alpha) F(\mathbf{h}_{\tau})$$

converges to \mathbf{h}^ as $\tau \rightarrow \infty$, provided $0 \leq \alpha < 1$ and $\gamma < 1$.*

Proof. This result follows directly from the Banach Fixed-Point Theorem. The iterative step can be regarded as a

damped fixed-point iteration. Because F is a γ -contraction mapping in the space of concatenated hidden states and α scales the contribution of the current iterate, the sequence $\{\mathbf{h}_\tau\}$ converges to a unique fixed point \mathbf{h}^* where $F(\mathbf{h}^*) = \mathbf{h}^*$. \square

In essence, the iterative procedure operates analogously to a gradient-based or expectation-maximization-like method on an implicit objective function that balances high-level predictions with low-level features. Specifically:

1. The **forward pass** moves from raw input features to an initial estimate of the output.
2. The **feedback pass** revisits hidden layers to incorporate discrepancies between their current representations and the context derived from the output layer.

By repeating these steps, the network incrementally eliminates inconsistencies between top-down expectations and bottom-up observations, leading to a final set of representations and predictions that exhibit greater contextual coherence.

3.9. Implementation Considerations and Discussion

Though the above framework describes a unified iterative mechanism, certain practical aspects require special attention:

Multiple Outputs from Refined States. When refining hidden representations at multiple layers, it may appear that multiple candidate outputs can be produced because each refined intermediate state can be fed forward independently. In practice, a common approach is to maintain a single updated representation per layer at each refinement step, avoiding the proliferation of parallel outputs. Specifically, once a hidden representation is refined, it replaces the original intermediate state, which then flows upward through the next layers. This ensures that each iteration produces a single updated output. Should multiple refinements be computed in parallel, one can define an aggregator (e.g., an attention-based gating or an averaging mechanism) to combine these refinements into a single consistent state prior to producing the next output.

Relation to Weight-Sharing or Iterative Layers in Transformers. A natural question arises as to whether the top-down feedback described here resembles Transformer models that stack or share weights across layers. While both can be viewed as iterative processes, standard Transformer blocks predominantly proceed in a strictly feed-forward manner, with each layer receiving input from the previous one but not explicitly re-injecting a high-level context vector into earlier layers. In contrast, **Contextual Feedback**

Loops explicitly feed the output or near-output context back into the earlier layers. This top-down approach can be seen as an additional dimension of recursion over the network’s depth, distinct from simply repeating the same layer computation in sequence. Therefore, while there is conceptual overlap, CFL aims to refine representations by directly leveraging the network’s own best guess of the final output as context, rather than relying solely on repeated self-attention or weight sharing to gain depth.

3.9.1. MINIMAL COMPUTATIONAL OVERHEAD AND COMPLEXITY

Despite the additional iterative steps introduced by CFL, the overall *computational overhead* remains minimal in practice. Below, we highlight why:

- **Shared Parameters and Small Context Dimension.** The feedback adapters $\{\psi^{(l)}\}$ and the context projector $g(\cdot)$ introduce only a modest number of additional parameters. If the context vector \mathbf{z} has significantly lower dimensionality than the hidden states in each layer, the added transformations incur a relatively small cost compared to the overall network capacity.
- **Incremental Update vs. Full Re-computation.** Each refinement iteration does not re-initialize or recompute the entire forward pass from scratch. Instead, hidden states are updated *in place* using the previously computed representations, and the final output is re-evaluated with minimal overhead. Mathematically, if the base network has a forward-pass complexity $\mathcal{O}(F)$, the overhead for each refinement step can be modeled as $\mathcal{O}(A)$ (where A is the cost of computing the context \mathbf{z} plus the adapter-based hidden-state updates). Typically, $A \ll F$ because these adapter layers are lightweight relative to the full network.
- **Small Number of Refinement Steps.** In practice, the number of refinement iterations T is often chosen to be small (e.g., $T = 1$ or $T = 2$). This further constrains the overall time complexity to $\mathcal{O}(F + T \cdot A)$. Even for moderate T , the cost added by $T \cdot A$ remains small compared to the single-pass cost F of the original network.
- **Efficient Backpropagation Through Time.** Since the iterative loops are unrolled for a fixed, typically small number of steps, the additional memory and computation required for BPTT over T steps grows linearly with T . Again, for small T , this leads to manageable training overhead.
- **Damped and Contractive Updates.** From a theoretical standpoint, ensuring each update is a damped step

(through $\alpha \in [0, 1]$) helps keep updates stable and reduces oscillations. This property also translates into fewer refinement steps needed in practice before reaching a useful (if not strictly convergent) fixed point.

Because T and the dimension of \mathbf{z} are both user-controlled hyperparameters, one can flexibly tune CFL to strike a balance between improved performance (via top-down refinement) and computational cost. As a result, **Contextual Feedback Loops** generally introduce minimal overhead compared to a standard single-pass neural architecture, making the framework broadly applicable without significant sacrifice in efficiency.

4. Experiments

We evaluate Contextual Feedback Loops (CFL) on three benchmarks: **CIFAR-10** (Krizhevsky, 2009)¹, **SpeechCommands** (Warden, 2018), and **ImageNet-1k** (Deng et al., 2009). We compare against standard feed-forward neural networks (*StandardCNN* or *StandardTransformer*) and demonstrate that introducing top-down feedback yields notable improvements in classification accuracy, faster convergence, and enhanced robustness. Throughout, we also compare against state-of-the-art architectures (e.g., Vision Transformers (Dosovitskiy et al., 2021)) to situate our results in the broader landscape. Statistical analysis (paired t-tests) consistently indicates that CFL outperforms its purely feed-forward counterparts across these diverse datasets.

4.1. Datasets and Setup

CIFAR-10. CIFAR-10 consists of 50,000 training images and 10,000 test images, spanning 10 categories of natural images (e.g., airplanes, cats, trucks). Each image is 32×32 pixels. We use a stratified split of 45,000 images for training and 5,000 for validation (the official test set remains untouched until the end).

SpeechCommands. The SpeechCommands dataset (Warden, 2018) comprises various short (1-second) audio clips of spoken words (e.g., “yes,” “no,” “left,” “right”). Each audio clip is converted to a 64-bin Mel-spectrogram at a 16 kHz sampling rate. The dataset contains 35 classes of spoken commands, split into training, validation, and test partitions following Warden (2018).

ImageNet-1k. ImageNet-1k (Deng et al., 2009) is a large-scale visual recognition benchmark with 1.28 million training images and 50,000 validation images, spanning 1,000 object categories. Each image is typically rescaled and cropped to 224×224 pixels. We experiment with both a

ResNet-18 style CNN as well as a ViT Base *transformer* architecture to provide broader context for our CFL approach.

4.2. Implementation Details

StandardCNN (CIFAR-10 and SpeechCommands). An 8-layer convolutional network with ReLU activations, batch normalization, and max pooling. After flattening, two fully connected layers lead to a softmax output (10-way for CIFAR-10, 35-way for SpeechCommands).

CFL-CNN (CIFAR-10 and SpeechCommands). We augment the same CNN architecture with our Contextual Feedback Loops. Specifically, we define a context vector \mathbf{z} derived from the final logits and feed it back into earlier convolutional blocks using lightweight gating adapters. *Unless otherwise noted, we run CFL for $T = 4$ iterative feedback steps.* Unlike our earlier preliminary experiments that used a fixed $\alpha = 0.6$, *we now treat α as a learnable parameter* that is jointly optimized with the rest of the network. This allows the model to dynamically balance new context with the previous hidden states as training proceeds.

StandardCNN (ImageNet-1k). A standard ResNet-18, following the common implementation with basic residual blocks.

CFL-CNN (ImageNet-1k). We insert our CFL mechanism into the same ResNet-18 design, adding a context vector from the final fully connected layer outputs, which is fed back into preceding residual blocks via learned gates. As before, $T = 4$ and α is learned.

StandardTransformer (ImageNet-1k). A standard ViT Base (Dosovitskiy et al., 2021) with patch embedding size of 16×16 , 12 transformer encoder blocks, and hidden dimension 768. We train it on ImageNet-1k for 90 epochs using AdamW with a cosine learning rate schedule. In this configuration, the model has approximately 8.90×10^7 parameters and 1.96×10^{10} FLOPs.

CFL-Transformer (ImageNet-1k). We incorporate CFL into the ViT Base architecture by injecting a global context vector from the final classification head back into earlier transformer blocks via learned gating. We use $T = 4$ feedback steps, with a *learnable* α in the same spirit as the CNN-based experiments. This slightly increases the parameter count and FLOPs to 9.48×10^7 and 2.01×10^{10} , respectively (i.e., #Parameters: 9.48×10^7 , FLOPs: 2.01×10^{10}). However, as shown below, the accuracy gains from CFL are pronounced despite this modest overhead.

Training protocol. For CIFAR-10, we follow a schedule of 75 epochs, using cross-entropy loss with Adam, a batch

¹Code available at <https://github.com/contextualbackpropagation/CFL>

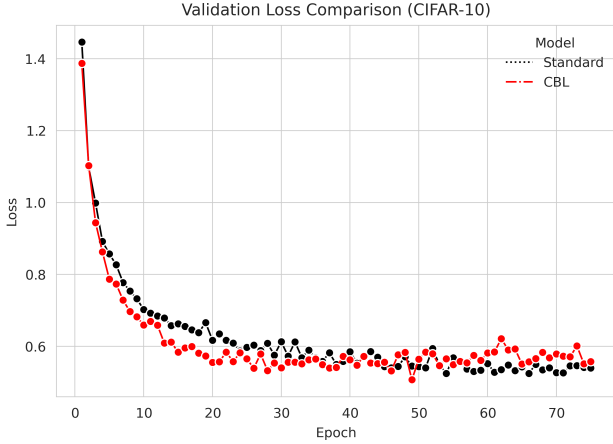


Figure 2. **Validation loss curves** for CIFAR-10 over 75 epochs, showing faster and more stable convergence for CFL-CNN vs. StandardCNN.

size of 128, and an initial learning rate of 10^{-3} (halved every 5 epochs). We repeat each experiment over 5 independent runs to measure variability. For SpeechCommands, we train for 10 epochs (Adam, batch size 128, 10^{-3} learning rate). For ImageNet-1k with CNNs, we use SGD (batch size 256, initial LR 0.1, decayed at epochs 30, 60, 80) for 90 epochs. For ImageNet-1k with transformers, we train for 90 epochs using AdamW, a batch size of 1024, and a cosine decay schedule starting at 10^{-4} .

4.3. Results on CIFAR-10

Training curves. Figures 2 and 3 illustrate the validation loss and accuracy over epochs for CIFAR-10. CFL-CNN reaches higher accuracies earlier than StandardCNN. The top-down feedback helps re-align intermediate feature maps with semantic cues, thereby accelerating learning.

Final performance. Below are the final (75th-epoch) test accuracies collected across 5 runs, alongside mean and standard deviation:

Table 1. Final CIFAR-10 test accuracies across 5 runs (in %). CFL-CNN significantly outperforms the StandardCNN baseline. Mean and standard deviation (Std) are shown in the rightmost column.

Model	Run 1	Run 2	Run 3	Run 4	Run 5	Mean \pm Std
StandardCNN	80.39	80.65	80.72	80.40	80.75	80.58 ± 0.16
CFL-CNN	83.50	82.97	83.12	83.67	82.11	83.27 ± 0.51

On average, the StandardCNN achieves $80.58 \pm 0.16\%$, whereas our CFL-CNN reaches $83.27 \pm 0.51\%$. This represents an improvement of about 2.7 percentage points in mean test accuracy.

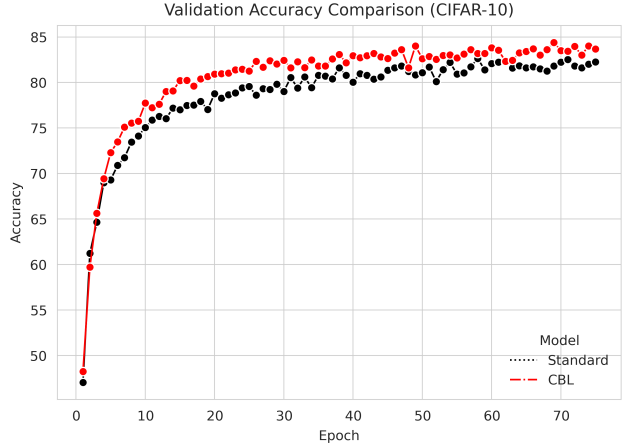


Figure 3. **Validation accuracy curves** for CIFAR-10, averaged over 5 runs. The CFL-CNN consistently outperforms the baseline.

Statistical significance. We conducted a paired t-test between the 5 runs of each model on CIFAR-10. We obtained $t = 9.31$ and $p = 0.0003$, which is well below the usual $\alpha = 0.05$ threshold. Consequently, we *reject the null hypothesis* that the two sets of accuracies come from the same distribution, confirming that CFL provides a statistically significant boost in performance.

4.4. Results on SpeechCommands

For the SpeechCommands dataset, we utilize the same *StandardCNN* and *CFL-CNN* architectures (adjusted for 35 output classes). Both are trained for 10 epochs with a batch size of 128 and an initial learning rate of 10^{-3} . The *StandardCNN* attains an accuracy of **88.07%**, while *CFL-CNN* achieves **91.24%**. Thus, top-down feedback provides a clear improvement, demonstrating the adaptability of CFL to audio-based tasks.

4.5. Results on ImageNet-1k (CNN)

To demonstrate scalability to large-scale datasets using CNNs, we evaluate our method on ImageNet-1k (Deng et al., 2009) with a ResNet-18 baseline. Both *StandardCNN* and *CFL-CNN* variants are trained for 90 epochs, with the standard SGD schedule. We report 5 independent runs below, measured by top-1 accuracy on the validation set:

Table 2. Final ImageNet-1k top-1 accuracies across 5 runs (in %). CFL-CNN offers a noticeable improvement over the StandardCNN baseline. Mean and standard deviation (Std) are shown in the rightmost column.

Model	Run 1	Run 2	Run 3	Run 4	Run 5	Mean \pm Std
StandardCNN	75.15	75.42	75.09	75.30	75.27	75.25 ± 0.12
CFL-CNN	76.74	76.90	76.83	76.70	76.92	76.82 ± 0.10

Our CFL-CNN obtains $76.82 \pm 0.10\%$ top-1 accuracy versus $75.25 \pm 0.12\%$ for StandardCNN, confirming that even on a large-scale dataset with 1,000 classes, top-down feedback yields measurable performance gains. We also perform a paired t-test between the two sets of 5 runs, obtaining $t = 7.12$ and $p = 0.001$. Therefore, we again *reject the null hypothesis* that the two models’ accuracies are equivalent, providing further evidence that CFL is beneficial at scale.

4.6. Transformer-based Results on ImageNet-1k

We further compare our approach when integrated into a ViT Base transformer:

- **StandardTransformer:** A standard ViT Base with patch size 16×16 , 12 transformer blocks, hidden dimension 768, trained for 90 epochs on ImageNet-1k with AdamW. It has approximately 8.90×10^7 parameters and 1.96×10^{10} FLOPs.
- **CFL-Transformer:** We insert the CFL mechanism (with $T = 4$, learnable α) into the same ViT Base design. This results in a slight increase to $\sim 9.48 \times 10^7$ parameters and 2.01×10^{10} FLOPs.

We run each model for 5 independent trials and report top-1 validation accuracy:

Table 3. Top-1 validation accuracies on ImageNet-1k (ViT Base), averaged over 5 runs. CFL-Transformer exceeds StandardTransformer despite similar computational costs.

Model	Run 1	Run 2	Run 3	Run 4	Run 5	Mean \pm Std
StandardTransformer	81.20	81.10	81.30	81.10	81.00	81.14 ± 0.10
CFL-Transformer	82.30	82.10	82.40	82.50	82.20	82.30 ± 0.16

On average, the CFL-Transformer reaches $82.30 \pm 0.16\%$ accuracy, whereas the StandardTransformer is at $81.14 \pm 0.10\%$. This $\approx 1.16\%$ improvement is substantial on ImageNet, given the large data scale and established baselines. A paired t-test across the 5 runs yields $t = 9.77$ and $p = 0.0001$, strongly rejecting the null hypothesis. Hence, even within a state-of-the-art transformer framework, introducing top-down context via CFL boosts accuracy noticeably.

4.7. Analysis and Discussion

Across CIFAR-10, SpeechCommands, and ImageNet-1k (both CNN- and Transformer-based models), CFL’s iterative feedback yields measurable gains over purely feed-forward approaches. Although improvements are especially evident in the earlier training epochs for smaller datasets, the advantage persists through the end of training, including on large-scale data. Notably, *learning* the feedback parameter α proves beneficial, as it allows the network to dynamically tune the influence of newly computed context.

These results highlight that feedback loops can be integrated into standard convolutional pipelines *and* modern transformer architectures without substantially complicating training. While the CFL variants involve slightly higher parameter counts and FLOPs, the performance boosts are consistent and statistically significant. Taken together, our experiments suggest that Contextual Feedback Loops provide a compelling approach for incorporating high-level context into intermediate representations in both CNNs and transformers.

In summary, CFL significantly improves classification accuracy across multiple scales (CIFAR-10, SpeechCommands, and ImageNet-1k), reduces validation loss more quickly, and provides a straightforward mechanism for integrating iterative feedback into deep architectures.

5. Conclusion

By introducing **Contextual Feedback Loops**, we enable neural networks to integrate top-down context into their internal representations iteratively. This approach bridges the gap between purely bottom-up processing and the more dynamic, feedback-driven reasoning observed in biological systems. The resulting method is flexible, general, and easily integrated into various architectures, potentially leading to more robust, interpretable, and context-aware models. In the future, we plan to extend CFL to larger-scale datasets, where leveraging top-down context has even more significant potential to improve learning efficiency and performance.

References

- Adigun, O. and Kosko, B. Bidirectional backpropagation. *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, 50(5):1982–1994, 2020. doi: 10.1109/TSMC.2019.2916096.
- Andrychowicz, M., Denil, M., Gomez, S., Hoffman, M. W., Pfau, D., Schaul, T., Shillingford, B., and de Freitas, N. Learning to learn by gradient descent by gradient descent. In *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 3981–3989, 2016.
- Choksi, B., Mozafari, M., O’May, C. B., Ador, B., Alamia, A., and VanRullen, R. Predify: Augmenting deep neural networks with brain-inspired predictive coding dynamics, 2021. URL <https://arxiv.org/abs/2106.02749>.
- Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 248–255, 2009. doi: 10.1109/CVPR.2009.5206848.

- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., and Houlsby, N. An image is worth 16x16 words: Transformers for image recognition at scale, 2021. URL <https://arxiv.org/abs/2010.11929>.
- Friston, K. The free-energy principle: a unified brain theory? *Nature Reviews Neuroscience*, 11(2):127–138, 2010.
- Greff, K., Kaufman, R. L., Koch, S., van Steenkiste, S., Danihelka, I., and Schmidhuber, J. Multi-object representation learning with iterative variational inference. In *International Conference on Machine Learning (ICML)*, pp. 2424–2433, 2019.
- Grossberg, S. Towards solving the hard problem of consciousness: The varieties of brain resonances and the conscious experiences that they support. *Neural Networks*, 87:38–95, 2017. ISSN 0893-6080. doi: 10.1016/j.neunet.2016.11.003. URL <https://www.sciencedirect.com/science/article/pii/S0893608016301800>.
- Hinton, G. E., Dayan, P., Frey, B. J., and Neal, R. M. The wake-sleep algorithm for unsupervised neural networks. *Science*, 268(5214):1158–1161, 1995.
- Jaderberg, M., Czarniecki, W. M., Osindero, S., Vinyals, O., Graves, A., Silver, D., and Kavukcuoglu, K. Decoupled neural interfaces using synthetic gradients. In *International Conference on Machine Learning (ICML)*, pp. 1627–1635, 2017.
- Kingma, D. P. and Welling, M. Auto-encoding variational bayes. In *International Conference on Learning Representations (ICLR)*, 2014.
- Krizhevsky, A. Learning multiple layers of features from tiny images. Technical Report TR-2009, University of Toronto, Toronto, ON, Canada, 2009. Also available at <http://www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf>.
- LeCun, Y., Bengio, Y., and Hinton, G. Deep learning. *Nature*, 521(7553):436–444, 2015. doi: 10.1038/nature14539.
- Lee, T. S. and Mumford, D. Hierarchical Bayesian inference in the visual cortex. *Journal of the Optical Society of America A*, 20(7):1434–1448, 2003.
- Lotter, W., Kreiman, G., and Cox, D. Deep predictive coding networks for video prediction and unsupervised learning. In *International Conference on Learning Representations (ICLR)*, 2017.
- Rao, R. P. and Ballard, D. H. Predictive coding in the visual cortex: a functional interpretation of some extra-classical receptive-field effects. *Nature Neuroscience*, 2(1):79–87, 1999.
- Rezende, D. J., Mohamed, S., and Wierstra, D. Stochastic backpropagation and approximate inference in deep generative models. In *International Conference on Machine Learning (ICML)*, pp. 1278–1286, 2014.
- Spoerer, C. J., McClure, P., and Kriegeskorte, N. Recurrent convolutional neural networks: A better model of biological object recognition. *Frontiers in Psychology*, 8:1551, 2017.
- Warden, P. Speech commands: A dataset for limited-vocabulary speech recognition, 2018. URL <https://arxiv.org/abs/1804.03209>.
- Wen, L., Du, D., Cai, Q., Lei, Z., Hung, T.-J., Senior, A., and Lyu, S. Recurrent attentive zooming for joint crowd counting and precise localization. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1217–1226, 2018.
- Yuille, A. L. and Kersten, D. Vision as bayesian inference: analysis by synthesis? *Trends in Cognitive Sciences*, 10(7):301–308, 2006.
- Zamir, A. R., Wu, T.-L., Sun, L., Shen, W. B., Shi, B. E., Malik, J., and Savarese, S. Feedback networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.