

Neuromodulatory Control Networks (NCNs): A Biologically Inspired Architecture for Dynamic LLM Processing

Author: Michael Morgan

Date: April 25, 2025

<https://github.com/Mmorgan-ML>

Abstract:

Large Language Models (LLMs) based on the Transformer architecture have achieved remarkable success, yet their core processing mechanisms remain largely static after training. While powerful, this static nature limits their ability to dynamically adapt their processing strategy based on nuanced contextual cues, task demands, or desired operational modes (e.g., shifting between exploration and exploitation). We propose **Neuromodulatory Control Networks (NCNs)**, a novel architectural modification inspired by the neuromodulatory systems in the vertebrate brain (e.g., those utilizing dopamine, acetylcholine, norepinephrine). NCNs are small, parallel networks that receive contextual input, summarizing the global state, task information, or external control signals, and compute dynamic "modulatory signals". These signals are broadcast to the main LLM to influence its computational properties during a forward pass, analogous to how neuromodulators alter neuronal gain, plasticity, and network states. Instead of merely routing information, NCNs aim to change *how* information is processed throughout the base model by modulating key components like attention mechanisms (e.g., temperature), layer gains, and activation functions. This paper introduces the NCN architecture, details its components and potential mechanisms, discusses its conceptual advantages for enhancing LLM adaptability, controllability, and efficiency, and provides an open-source PyTorch implementation to facilitate community exploration and future empirical validation.

1. Introduction

Large Language Models (LLMs), predominantly based on the Transformer architecture [22], have revolutionized natural language processing, demonstrating impressive capabilities in generation, understanding, and reasoning. Foundational models like GPT [19, 2] and BERT [5] exemplify this success. However, despite their scale and power, the fundamental processing pathway within these models is typically fixed after pre-training or fine-tuning. Each input token, regardless of the broader context or specific task requirements, undergoes a largely identical sequence of computations dictated by static model parameters. This inherent staticity contrasts sharply with the dynamic adaptability observed in biological brains, which leverage neuromodulatory systems to flexibly alter processing modes based on internal states and external cues [6]. Current LLMs often struggle with fine-grained contextual adaptation within a single processing sequence, exhibit interference in multi-task scenarios, and offer limited avenues for precise behavioral control beyond modifications to the input prompt [2, 14].

Inspired by the principles of biological neuromodulation, we propose **Neuromodulatory Control Networks (NCNs)** as an architectural enhancement for LLMs. Neuromodulatory systems, originating from small nuclei but projecting broadly, release neurotransmitters like dopamine (DA), acetylcholine (ACh), norepinephrine (NE), and serotonin (5-HT). These modulators orchestrate global changes in brain state, influencing attention, learning rates, arousal, memory encoding versus retrieval, and the balance between exploration and exploitation [6, 20]. They dynamically alter how neural circuits process information, rather than solely carrying the primary information content themselves.

Our proposed NCN architecture mirrors this principle. It comprises a standard base LLM (e.g., a Transformer [22]) augmented with one or more small, parallel NCN modules. These NCNs receive input summarizing the current context, which could include:

- Aggregate representations of the input sequence (e.g., pooled hidden states).
- Task identifiers or instruction embeddings (for multi-task settings).
- Statistics derived from the LLM's internal activations.
- Explicit external control signals specifying desired behavior (e.g., style, focus).

Based on this context, the NCN computes a set of "modulatory signals"—low-dimensional vectors or scalars. These signals are then broadcast throughout the base LLM to dynamically adjust the functioning of its core components during the forward pass. Key proposed modulation mechanisms include:

- **Modulating Attention:** Adjusting attention softmax temperature to control focus (sharp vs. broad) or dynamically gating/mixing attention head outputs.

- **Modulating Layer Gains:** Scaling the output of specific layers or components (e.g., attention blocks, feed-forward networks) to increase or decrease their influence.
- **Modulating Activation Functions:** Altering the shape or gating parameters of activation functions (e.g., within FFNs like SwiGLU [20] or GELU [10]).
- *(Potentially)* **Modulating Learning Dynamics:** During training or fine-tuning, using NCN outputs to set context-dependent learning rates or regularization strengths.

This approach fundamentally differs from existing methods for conditional computation. Unlike Mixture-of-Experts (MoE) which routes tokens to distinct, static experts [22, 7, 15], NCNs aim to modulate the processing within *shared* components. Compared to Adapters [12, 18, 9] or Hypernetworks [8, 16] that typically generate *static*, task-specific parameters or modules, NCNs provide *dynamic* modulation signals that can potentially vary throughout the processing of a single input sequence based on evolving context. While incorporating gating mechanisms similar to those in LSTMs [11], GRUs [3], or FiLM layers [16], NCNs compute these gates based on a potentially more global context captured by the dedicated NCN module, rather than solely local state. Some prior work has also explored dynamically adjusting representations based on context, but NCNs propose a broader, neuromodulation-inspired framework for controlling multiple aspects of the LLM processing itself.

We hypothesize that NCNs offer several potential benefits, including enhanced contextual adaptability (e.g., switching between factual recall and creative generation modes), improved multi-task learning through dynamic parameter sharing, more integrated controllability via exposed NCN inputs (potentially interfacing with methods like PPLM [4], CTRL [13], or RLHF [23]), and potential for dynamic resource allocation.

The primary contributions of this work are:

1. The proposal of the Neuromodulatory Control Network (NCN) architecture, grounded in principles of biological neuromodulation [6, 20], for dynamically adapting LLM processing.
2. The definition and formulation of specific mechanisms by which NCN-generated signals can modulate core Transformer components (attention, FFNs [using activations like GELU [10] or SwiGLU [20]], gains, Layer Normalization [1]).
3. An open-source PyTorch implementation of the NCN architecture to serve as a foundation for future research and experimentation by the community.
4. A conceptual analysis positioning NCNs relative to existing dynamic and conditional computation techniques in LLMs [22, 12, 8, 17].

This paper details the NCN architecture, its biological motivation, potential mechanisms, and implementation considerations. While empirical validation remains future work, we believe NCNs

represent a promising direction for developing more flexible, context-aware, and controllable language models. We hope this proposal and the accompanying codebase will stimulate further investigation into neuromodulation-inspired AI architectures.

The remainder of this paper is structured as follows: Section 2 discusses biological neuromodulation and related work in dynamic LLM architectures. Section 3 details the proposed NCN architecture and modulation mechanisms. Section 4 provides implementation notes based on our codebase. Section 5 outlines potential benefits and challenges. Section 6 compares NCNs to existing approaches. Section 7 suggests directions for future research, and Section 8 concludes.

2. Background and Related Work

To contextualize the Neuromodulatory Control Network (NCN) proposal, we first review the principles of biological neuromodulation that inspire its design. We then discuss existing approaches within machine learning, particularly for Large Language Models (LLMs), that aim to achieve dynamic adaptation, conditioning, or control, highlighting the conceptual distinctions from NCNs.

2.1 Biological Neuromodulation: Dynamic Brain State Control

Neural computation in the brain relies on both fast, targeted synaptic transmission (primarily mediated by glutamate and GABA) and slower, broader neuromodulation. Neuromodulatory systems typically involve relatively small clusters of neurons, often located in the brainstem or basal forebrain, that project diffusely to vast areas of the central nervous system, including the cortex, hippocampus, thalamus, and basal ganglia [6]. Instead of directly transmitting specific informational content, neuromodulators like dopamine (DA), acetylcholine (ACh), norepinephrine (NE), and serotonin (5-HT) act more like "broadcast signals", altering the computational properties and response dynamics of the circuits they innervate [6, 20].

Key neuromodulators and their associated functions provide inspiration for NCN mechanisms:

- **Dopamine (DA):** Strongly associated with reward processing, motivation, and reinforcement learning [19]. Phasic DA signals often encode "reward prediction error"—the difference between expected and received rewards—serving as a crucial learning signal, particularly in the basal ganglia and prefrontal cortex, to reinforce beneficial actions or state transitions. DA also plays roles in modulating working memory gating and attentional control.
- **Acetylcholine (ACh):** Critical for attention, learning, and memory formation. ACh release, often triggered by novel or behaviorally salient stimuli, enhances sensory processing, increases neuronal responsiveness (gain modulation), and promotes synaptic plasticity,

particularly within the hippocampus and cortex [6]. It is thought to shift network states to favor encoding new information.

- **Norepinephrine (NE) / Noradrenaline (NA):** Involved in arousal, alertness, vigilance, and the "fight-or-flight" response. NE release, often during stress or salient events, can sharpen attentional focus, enhance signal-to-noise ratios in neural processing, and modulate memory consolidation. It is linked to shifts between exploratory behavior (seeking new information) and exploitative behavior (using known strategies) [6].
- **Serotonin (5-HT):** Possesses perhaps the most diverse and widespread effects, implicated in mood regulation, sleep-wake cycles, patience, and social behavior. Its cognitive roles are complex, potentially influencing behavioral inhibition, temporal discounting, and modulating the impact of negative feedback or punishment, thereby affecting the overall stability and flexibility of network states [6].

From these biological systems, we extract several key principles relevant for designing adaptive artificial networks:

1. **Context-Dependency:** Neuromodulatory effects are not fixed; they depend heavily on the current behavioral state, task demands, and environmental context. The same neuron or circuit can respond differently under varying neuromodulatory conditions.
2. **Diffuse Influence:** Modulators often affect large populations of neurons simultaneously, changing the overall "tone" or processing mode of entire brain regions or distributed networks.
3. **Gain Modulation:** A common effect is altering the gain of neuronal responses (the slope of the input-output function), making neurons more or less sensitive or responsive to their inputs (akin to ACh/NE effects on attention/arousal).
4. **Plasticity Modulation:** Neuromodulators act as critical "gates" for synaptic plasticity, determining *when* and *how strongly* connections should change based on neural activity and outcomes (akin to DA's role in learning).
5. **State Switching:** They facilitate transitions between distinct functional states (e.g., exploration vs. exploitation, focused attention vs. broad monitoring, learning vs. recall).

NCNs aim to capture these principles computationally, providing a mechanism for LLMs to dynamically adjust their processing based on context, analogous to how these biological systems shape brain function.

2.2 Related Work in Dynamic and Conditional LLMs

Several existing techniques aim to make LLMs more adaptive, controllable, or efficient. NCNs share the goal of context-dependent processing but differ mechanistically.

- **Prompting and In-Context Learning:** Prepending specific instructions or examples (prompts) to the input is the standard way to condition the behavior of large pre-trained models [2]. This relies on the model's ability to recognize patterns in the prompt and adapt its subsequent generation through its standard forward pass, without any architectural changes or explicit dynamic modulation during processing. While effective, especially at scale, control is implicit and depends heavily on model size and prompt formulation. NCNs propose an explicit, internal architectural mechanism for adaptation based on context signals derived potentially beyond the immediate prompt.
- **Mixture-of-Experts (MoE):** MoE models increase capacity by employing multiple expert subnetworks (typically FFN layers) and using a gating network to route each input token to a small subset of experts (often just one or two) [22, 7, 15]. This allows for scaling model size while keeping computational cost per token roughly constant. NCNs differ fundamentally: instead of routing tokens to different *static* experts, NCNs compute signals that modulate the behavior of *shared* processing components experienced by all tokens. While an NCN could potentially modulate the MoE gating network itself, its core proposed function is to alter the properties of the underlying pathways, not just select among them.
- **Adapters and Parameter-Efficient Fine-Tuning (PEFT):** Adapters involve inserting small, trainable modules into the layers of a fixed pre-trained model [12, 18]. Only the adapter parameters are trained for downstream tasks, allowing efficient specialization. Various PEFT methods exist [9]. These approaches typically adapt the model *statically* for a given task or dataset. NCNs, in contrast, aim for *dynamic* adaptation within a single forward pass based on the immediate context, modulating the behavior of the (potentially fixed) base model parameters rather than inserting separate static modules. NCNs could potentially be combined with adapters, perhaps by having NCNs modulate adapter behavior or having adapters fine-tune an NCN-augmented base model.
- **Hypernetworks:** Hypernetworks are networks that generate the weights (or parameters) for another primary network [8]. In the context of LLMs, they have been used, for example, to generate task-specific adapter parameters [15]. While context-dependent, hypernetworks typically generate a *static* set of weights for the primary network based on a given task or input instance *before* the main processing begins. NCNs differ by generating *dynamic modulation signals* that continuously influence the primary network's processing throughout its forward pass, potentially changing based on evolving internal state or external inputs.
- **Conditional Parameters and Gating:** The idea of dynamically computed parameters influencing network behavior has precedents. Recurrent Neural Networks like LSTMs [11] and GRUs [3] use internal gating mechanisms computed from the local input and hidden state to control information flow. FiLM (Feature-wise Linear Modulation) layers use conditioning information to compute per-feature scaling (gain) and shifting (bias) parameters,

often applied in vision-language tasks [16]. Other related work also explores dynamically constructing representations based on context. NCNs build on this concept but frame it within the broader context of neuromodulation, typically computing modulation signals based on a more global context (pooled states, task IDs, external controls) via a dedicated NCN module and applying these signals potentially to a wider range of targets (attention temperature, activation shapes) beyond simple gain/bias or LSTM/GRU-style local gating.

- **Controllable Generation:** Techniques exist to steer LLM outputs towards desired attributes, such as specific topics, styles, or toxicity levels. Methods include using control codes prepended to the input [13], modifying activations during generation [4], or employing reinforcement learning from human feedback (RLHF) [23]. NCNs offer a potential architectural substrate for such control; external signals representing desired attributes could be fed directly into the NCN, allowing it to generate corresponding modulation signals that shape the LLM's internal processing to achieve the desired output characteristics in a potentially more integrated manner.

In summary, while sharing goals with existing methods, the proposed NCN architecture offers a distinct, biologically inspired approach. It focuses on dynamically altering the *processing mode* of shared LLM components using broadcast signals computed by a dedicated contextual network, aiming for a flexible and integrated form of adaptation and control.

3. Neuromodulatory Control Network Architecture

We propose augmenting a standard Large Language Model (LLM), typically a Transformer [22], with one or more parallel Neuromodulatory Control Networks (NCNs). The NCN is designed to be significantly smaller and computationally cheaper than the main LLM. Its role is to dynamically compute context-dependent modulation signals that influence the processing within the main LLM layers during a single forward pass.

3.1 Overall Framework

The core architecture consists of two main components operating in parallel:

1. **Base LLM:** A standard Transformer network composed of embedding layers, multiple Transformer blocks (each containing self-attention and feed-forward sub-layers with residual connections and normalization [1]), and an output prediction head.
2. **Neuromodulatory Control Network (NCN):** A smaller neural network that receives contextual input derived from the main LLM's state or external sources. It outputs a set of modulation signals.

The NCN computes its modulation signals based on its input, and these signals are then used within each layer of the base LLM to modify specific computational steps, such as attention scoring or sub-layer outputs. This modulation aims to dynamically alter *how* information is processed by the fixed weights of the base LLM.

3.2 NCN Input Formulation

The NCN requires access to information that signals the current context or desired processing state. The specific input representation can be varied depending on the application. Based on our proposal and implementation (`model.py`), potential inputs include:

- **Global Sequence Representation:** An aggregate representation of the input sequence, computed from the initial embedding layer or an early Transformer layer. A common approach, implemented in our codebase, is to use masked average pooling over the token embeddings, ensuring padding tokens are excluded:

```
ControlInput = Pool(Embeddings * AttentionMask)
```

where `Pool` can be mean or max pooling. This provides the NCN with a summary of the entire input sequence content.
- **Task Information:** For multi-task scenarios, an embedding representing the current task ID or natural language instructions describing the task can be concatenated with or used as input to the NCN. This allows the NCN to tailor modulation signals specifically for the task at hand.
- **Internal State Metrics:** Statistics computed from the main LLM's activations during processing could serve as input. Examples include average activation magnitudes in certain layers, the entropy or variance of attention distributions (indicating focus vs. diffusion), uncertainty estimates from the prediction head, or gradient norms during training. These inputs would allow the NCN to react to the LLM's internal processing dynamics.
- **External Control Signals:** Explicit signals provided by a user or a controlling system can be fed into the NCN to steer LLM behavior. Examples include desired sentiment, style (e.g., formal vs. creative), topic focus, or safety level parameters.

Our reference implementation primarily utilizes the global sequence representation (masked average pooling of embeddings) as the NCN input (see `model.py`, `forward` method). The input dimension expected by the NCN (`ncn_input_dim`) is configurable via `NCNConfig` and defaults to the base LLM's hidden dimension (`d_model`).

3.3 NCN Architecture

The NCN module itself should be lightweight. Our reference implementation (`ncn.py`) utilizes a simple Multi-Layer Perceptron (MLP) as configured by `NCNConfig`:

```
RawSignals = MLP(ControlInput)
```

```
MLP(x) = Linear_2(Activation(Linear_1(x)))
```

where `Linear_1` maps from `ncn_input_dim` to `ncn_hidden_dim`, `Activation` is a non-linearity (e.g., ReLU, GeLU, specified by `ncn_activation_fn` in `NCNConfig`), and `Linear_2` maps from `ncn_hidden_dim` to `num_mod_signals` (the total number of raw outputs).

Alternative NCN architectures could include Recurrent Neural Networks (RNNs, like LSTMs [11] or GRUs [3]) if temporal integration of control inputs across a sequence or dialogue is needed, or even a very small Transformer.

3.4 Modulation Signals and Transformations

The raw output vector from the NCN (`RawSignals`, shape `(batch_size, num_mod_signals)`) is typically not used directly. Instead, different slices or elements of this vector are transformed to produce specific modulation signals constrained to appropriate ranges for their intended function. The mapping from raw outputs to named signals and their transformations is handled within the `NeuromodulatoryControlNetwork` module (`ncn.py`), guided by the `modulation_signal_names` list in `NCNConfig`.

Our implementation defines transformations for signals intended to mimic biological effects:

- **Gain Signals** (To scale layer outputs, often desired to be centered around 1.0. A sigmoid function followed by an offset achieves this:
`gain = sigmoid(raw_slice) + 0.5` (Range: [0.5, 1.5])
- **Attention Temperature** (To scale attention logits; must be positive. A softplus function ensures positivity, with a small epsilon for numerical stability:
`attention_temp = softplus(raw_slice) + epsilon` (Range: (epsilon, +inf), e.g., epsilon=1e-4)
- **Gating Signals** (To multiplicatively gate activations or outputs, typically in the range [0, 1]. A sigmoid function is suitable:
`ffn_gate = sigmoid(raw_slice)` (Range: [0, 1])

These transformations ensure that the modulation signals operate within sensible bounds, preventing instability (e.g., negative temperature, exploding gains) and aligning with their intended purpose. The processed signals are returned as a dictionary mapping signal names (e.g., "gain") to their corresponding tensors (e.g., shape `(batch_size, 1)`).

3.5 Modulation Targets and Mechanisms

The crucial part is how the processed modulation signals dynamically influence the base LLM's computation. These mechanisms must be differentiable to allow for end-to-end training. Our implementation (`ModulatedTransformerLayer.py`, `attention.py`, `feedforward.py`) focuses on modulating standard Transformer components within a Pre-LayerNorm structure:

- **Modulating Attention Temperature:** The `attention_temp` signal is used to divide the pre-softmax attention scores ($QK^T / \sqrt{d_k}$), implemented within our custom `MultiHeadAttention` module (`attention.py`):

```
ScaledScores = Scores / attention_temp.clamp(min=epsilon)
```

Lower temperatures lead to sharper, more focused attention distributions, while higher temperatures broaden the distribution, mimicking attentional effects of ACh/NE.
- **Modulating Layer Gains:** The `gain` signal is used to multiplicatively scale the output of sub-layers before the residual connection. In our implementation (`transformer_layer.py`), the *same* gain signal modulates the output of both the self-attention block and the feed-forward block within a layer:

```
AttnOutput_Modulated = AttnOutput * gain  
FFNOutput_Modulated = FFNOutput * gain
```

This increases or decreases the overall influence of the layer's computations based on context, analogous to ACh/NE gain modulation.
- **Modulating FFN Output:** The `ffn_gate` signal is passed to the `PositionwiseFeedForward` module (`feedforward.py`) and, in the current implementation, applied multiplicatively to its output:

```
FinalFFNOutput = BaseFFNOutput * ffn_gate (Potentially combined with gain)
```

This allows the NCN to control the contribution of the FFN based on context. Alternatively, if using architectures like SwiGLU [20], the NCN could modulate the internal gate activation directly.

While not implemented in the current codebase, other potential targets include modulating activation function parameters (e.g., slope of LeakyReLU), adding learned biases to Query/Key/Value projections in attention, dynamically mixing attention head outputs, or modulating learning rates/regularization during training.

3.6 Scope and Frequency of Modulation

A design choice involves how often modulation signals are computed and how broadly they apply. Our reference implementation adopts a simple and efficient approach:

- **Frequency:** The NCN is computed *once per input sequence*, based on the initial NCN input (e.g., pooled embeddings).
- **Scope:** The *same dictionary* of modulation signals is passed to *all* Transformer layers in the base LLM.

This provides a global, sequence-level modulation. More complex variants could involve:

- Recomputing NCN signals at each layer, potentially using intermediate LLM states as input.
- Computing signals on a per-token basis (computationally expensive).
- Generating layer-specific signals (e.g., different gain values for early vs. late layers).

The current approach balances simplicity and computational cost with providing a global contextual adaptation signal.

3.7 Interaction with Normalization Layers

The placement of modulation relative to normalization layers like Layer Normalization [1] is important. Our implementation uses a standard Pre-Norm structure (`output = residual + Dropout(SubLayer(LayerNorm(input)))`). Modulation (e.g., gain) is applied to the output of the sub-layer (`AttnOutput, FFNOutput`) *before* the dropout and residual addition. This means the subsequent `LayerNorm` at the input of the *next* sub-layer will potentially rescale the activations, interacting with the applied gain. Studying this interaction and potentially exploring alternative placements (e.g., modulating *after* `LayerNorm`, or developing adapted normalization schemes) is an area for future investigation.

4. Implementation Details

To facilitate understanding, experimentation, and adoption of the proposed architecture, we provide an open-source reference implementation using PyTorch. The codebase, structured within the `ncn_project` directory, includes the core NCN model architecture and a custom-developed tokenizer implementation necessary for practical application. This section describes the overall structure and highlights key aspects of the implementation. The full codebase is available at [Link to your GitHub Repository - **Insert URL Here**] and is licensed under the Creative Commons Attribution-ShareAlike 4.0 International License, with the exception of the GPT-2 tokenizer data files whose original license may differ.

4.1 Code Repository Structure

The `ncn_project` repository is organized into three main subdirectories:

1. `ncn_architecture/`: Contains the core PyTorch implementation of the NCN-modulated LLM, representing the primary architectural contribution detailed in this paper.
2. `gpt2_tokenizer_files/`: Holds standard data files (`vocab.json`, `merges.txt`, `tokenizer_config.json`, etc.) defining a GPT-2 style Byte-Pair Encoding (BPE) tokenizer. These files originate from the standard GPT-2 tokenizer configuration.
3. `tokenizer/`: Contains a **custom Python implementation**, developed by the author as part of this project, for loading and utilizing the BPE tokenizer defined by the configuration files in `gpt2_tokenizer_files/`. This includes modules for the tokenizer model logic, normalizers, pre-tokenizers, and processors necessary to convert raw text into model inputs. This custom tokenizer code is released under the same CC BY-SA 4.0 license as the `ncn_architecture` module.

While this paper concentrates on the novel model architecture within `ncn_architecture/`, the inclusion of the original GPT-2 tokenizer data files and the custom Python tokenizer implementation provides a complete, self-contained package for users wishing to experiment with the model.

4.2 Core Architecture Implementation (`ncn_architecture/`)

The `ncn_architecture/` directory contains the Python modules directly implementing the NCN-LLM:

- `config.py`: Defines the `NCNConfig` dataclass, centralizing hyperparameters for the base Transformer and NCN module.
- `ncn.py`: Implements the `NeuromodulatoryControlNetwork` class (MLP-based) and handles the transformation of raw outputs into named, scaled modulation signals.
- `attention.py`: Provides the custom `MultiHeadAttention` class capable of applying `attention_temp` modulation.
- `feedforward.py`: Contains the `PositionwiseFeedForward` class with optional output gating/scaling via `ffn_gate`.
- `transformer_layer.py`: Defines the `ModulatedTransformerLayer`, integrating the modulated components within a Pre-LayerNorm structure.
- `model.py`: Implements the main `ModulatedLLM` class, orchestrating all components, computing NCN input, and preparing masks.

Key implementation choices within this core module, as detailed in Section 3, include the MLP-based NCN, masked average pooling for NCN input, specific modulation mechanisms for attention temperature and layer/FFN gain/gating, Pre-LayerNorm structure, and standard weight initialization and tying practices.

4.3 Custom Tokenizer Implementation (tokenizer/ & gpt2_tokenizer_files/)

To enable practical use, the repository includes components to tokenize raw text:

- The files in `gpt2_tokenizer_files/` provide the necessary vocabulary and merge rules data for a standard GPT-2 BPE tokenizer.
- The code within the `tokenizer/` directory is a **custom implementation** developed for this project. It provides the Python classes and logic required to load the data from `gpt2_tokenizer_files/` and perform text tokenization (including normalization, pre-tokenization, BPE model application, and processing) to produce `input_ids` and `attention_mask` tensors compatible with the `ModulatedLLM`.

This custom tokenizer allows users to work with the `ModulatedLLM` without external dependencies for this specific GPT-2 style tokenization process and is provided under the CC BY-SA 4.0 license.

4.4 Usage Example (Conceptual)

The following conceptual example illustrates how the components provided in the repository might be used together, highlighting the use of the custom tokenizer:

```
# Example Conceptual Usage Reflecting Repository Structure
from ncn_architecture import NCNConfig, ModulatedLLM
# Import the custom tokenizer class from the 'tokenizer' module
from tokenizer import NCNTokenizer # Assuming NCNTokenizer is the main class name
import torch

# --- Tokenizer Loading ---
# Instantiate the custom tokenizer, potentially pointing it to the data files
tokenizer_data_dir = "path/to/ncn_project/gpt2_tokenizer_files"
tokenizer = NCNTokenizer(data_dir=tokenizer_data_dir) # Conceptual initialization

# --- Model Configuration and Instantiation ---
# Configure the model, ensuring vocab size matches the loaded tokenizer
config = NCNConfig(
    vocab_size=tokenizer.get_vocab_size(), # Get vocab size from tokenizer instance
    d_model=768,
    nhead=12,
    num_layers=12,
)
# Instantiate the model
model = ModulatedLLM(config)

# --- Input Preparation ---
text = "Example input text to be processed."
```

```

# Use the custom tokenizer's method to encode text
inputs = tokenizer.encode(text) # Conceptual encoding method
input_ids = inputs["input_ids"].unsqueeze(0) # Add batch dimension
attention_mask = inputs["attention_mask"].unsqueeze(0) # Add batch dimension

# --- Forward Pass ---
# Model requires input_ids and optional attention_mask
logits = model(input_ids=input_ids, attention_mask=attention_mask)

# logits shape: (1, sequence_length, config.vocab_size)

```

This reference implementation, featuring the core NCN architecture and the accompanying custom-developed tokenizer components, provides a robust foundation for future research, empirical evaluation, and community development based on the NCN concept.

5. Potential Benefits and Challenges

The proposed Neuromodulatory Control Network (NCN) architecture, while awaiting empirical validation, offers several potential advantages over static LLM architectures, alongside inherent challenges that warrant consideration for future research and development.

5.1 Potential Benefits

We hypothesize that incorporating NCNs could lead to significant improvements in LLM capabilities:

1. **Enhanced Contextual Adaptability:** This is the primary expected benefit. By dynamically modulating processing based on context signals, NCN-augmented LLMs could adjust their computational strategy on the fly. For instance, during factual question answering, the NCN might learn to lower attention temperature [Section 3.5] to sharpen focus on relevant details and potentially increase gain in layers associated with knowledge retrieval. Conversely, for creative writing tasks, it might broaden attention (higher temperature) and perhaps increase gain in later, more abstract layers to foster diverse and novel continuations. In dialogue systems, modulation could shift based on detected user intent, emotional tone, or dialogue history length, leading to more nuanced and appropriate responses.
2. **Improved Multi-Task Learning:** Standard multi-task learning often relies solely on parameter sharing or static adaptation methods like adapters [12, 18]. NCNs offer a mechanism for *dynamic* parameter sharing, where the shared base model parameters are configured differently for each task via the NCN's modulation signals (potentially derived from task embeddings fed to the NCN [Section 3.2]). This could lead to less interference between tasks and improved overall performance compared to static approaches, as the model can

adopt task-specific processing modes using the same underlying weights.

3. **Integrated Control and Steerability:** Exposing control signals as direct inputs to the NCN [Section 3.2] provides a potentially more integrated way to steer LLM behavior compared to solely relying on prompt engineering [2]. Imagine adjusting sliders for "creativity," "formality," or "safety level" that directly influence NCN inputs, which in turn shape the internal processing dynamics (e.g., attention temperature, layer gains) via modulation signals to achieve the desired output style or constraints. This contrasts with methods like PPLM [4] or CTRL [13] by embedding control within the core processing loop via the NCN architecture.
4. **Dynamic Resource Allocation:** Modulation signals, particularly gain or gating factors [Section 3.4, 3.5], could potentially learn to reduce or effectively nullify the contribution of specific components (e.g., attention heads, FFN layers, or even entire Transformer layers) when they are not deemed necessary for the current context or task. If a gain signal approaches zero, the computational contribution of that component is minimized. This opens possibilities for conditional computation, potentially leading to significant efficiency gains during inference by dynamically pruning computation paths based on NCN outputs, similar in spirit but different in mechanism to MoE routing [7, 22].
5. **More Biologically Plausible Learning:** If NCNs are configured to modulate learning dynamics (e.g., layer-specific learning rates or regularization strengths) during training [Section 3.5], this could enable dynamic, context-sensitive learning rules that align better with biological principles [6, 20] than fixed schedules. This might be particularly relevant for continual learning scenarios, where the NCN could adjust plasticity based on task novelty or perceived interference.
6. **Adaptive Exploration/Exploitation Balance:** Factors like attention temperature or potentially direct modulation of output sampling parameters could be controlled by the NCN. This would allow the model to dynamically shift between generating diverse, exploratory outputs (e.g., higher temperature, broader attention) and converging on high-probability, exploitative outputs (e.g., lower temperature, focused attention), adapting its generation strategy based on task requirements or explicit control signals.

5.2 Challenges and Considerations

Despite the potential benefits, the NCN architecture presents several challenges that need to be addressed:

1. **Stability:** Unconstrained modulation signals pose a significant risk to network stability. Gains that grow too large could cause activation explosion, while temperatures approaching zero could lead to numerical issues and overly peaky distributions. Careful design of the NCN's output transformation functions (e.g., using sigmoid, tanh, softplus with appropriate

scaling/offsets, as in Section 3.4) is crucial. Furthermore, the interaction between dynamic modulation and existing normalization layers (e.g., LayerNorm [1]) requires careful study [Section 3.7], as normalization might dampen or counteract the intended modulation effects, potentially necessitating adapted normalization schemes or specific placement of modulation relative to normalization.

2. **Training Dynamics:** Training the NCN end-to-end with the main LLM might be challenging. Gradients flowing back from the task loss through the modulated components into the NCN could become weak or noisy, making learning difficult. The NCN might converge to trivial solutions (e.g., always outputting default "neutral" modulation signals) unless the training task explicitly demands dynamic adaptation. Careful tuning of relative learning rates for the NCN versus the base LLM might be necessary. More complex training regimes, such as using auxiliary losses designed to encourage meaningful modulation or employing reinforcement learning [23] (treating the NCN as an agent choosing modulation actions), might be required but would add complexity.
3. **Interpretability:** While inspired by specific neuromodulators [6, 20], the signals learned by the NCN might not map cleanly onto distinct DA/ACh/NE/5-HT functions. Understanding *why* the NCN chooses certain modulation patterns in response to specific inputs will require dedicated analysis tools (e.g., correlating NCN signal values with input features, task performance, or internal LLM states), posing a significant interpretability challenge.
4. **Computational Overhead:** Although the NCN itself is designed to be small, computing its output and applying the modulation signals at each layer adds computational cost (FLOPs and potentially latency) to every forward pass. This overhead must be weighed against the performance gains or efficiency benefits (e.g., from dynamic resource allocation).
5. **NCN Input Design:** The effectiveness of the NCN critically depends on receiving the right contextual information [Section 3.2]. Determining the optimal NCN inputs (e.g., global vs. local state, inclusion of specific task info or internal metrics) is likely to require empirical investigation and may be task-dependent.
6. **Modulation Scope and Granularity:** Deciding whether modulation signals should be global (same signal for all layers), layer-specific, module-specific (e.g., separate gain for attention vs. FFN), or even more granular (e.g., per-head) involves trade-offs [Section 3.6]. Global modulation is simplest but least flexible; highly granular modulation offers flexibility but increases complexity and the risk of overfitting. Layer-specific or module-specific modulation seems a reasonable compromise but requires careful definition of the NCN output structure and application points.
7. **Initialization Strategy:** Proper initialization of the NCN is likely important for stable training [Section 4.2]. Initializing the NCN such that it initially outputs "neutral" modulation signals (e.g., gains/gates of 1.0, temperatures of 1.0) could allow the main LLM to establish baseline

performance before significant modulation occurs, preventing destabilization early in training.

Addressing these challenges will be crucial for realizing the full potential of the NCN architecture. Future work should focus on systematically investigating these aspects through empirical studies.

6. Comparison to Existing Approaches

Neuromodulatory Control Networks (NCNs) aim to enhance LLM adaptability and control through dynamic, context-dependent modulation of core processing components. While sharing goals with several existing techniques, the proposed NCN architecture differs significantly in its mechanism and scope. This section contrasts NCNs with prominent related approaches.

- **Conditioning via Prompts:** Standard prompting modifies the input sequence to guide LLM behavior [2]. Control is achieved implicitly through the model's learned ability to condition on preceding tokens. NCNs, conversely, introduce an *explicit internal mechanism* for adaptation. The NCN directly computes modulation signals based on broader context (potentially including, but not limited to, the prompt), which then actively alter the internal processing dynamics of the fixed base model parameters, rather than relying solely on the model's implicit interpretation of input patterns.
- **Mixture-of-Experts (MoE):** MoE models achieve conditional computation by routing tokens to different specialized "expert" subnetworks (typically FFNs) using a gating mechanism [22, 7, 15]. The core idea is selecting *which* expert processes a token. NCNs operate differently: they typically aim to modulate the computational properties of *shared* pathways experienced by all tokens within a layer (though NCNs could potentially modulate MoE gating). The goal is to change *how* information is processed by the base components (e.g., altering attention focus, gain), rather than selecting *which* static component performs the processing.
- **Adapters and Parameter-Efficient Fine-Tuning (PEFT):** Adapters [12, 18] and related PEFT methods [9] insert or modify a small number of parameters within a largely frozen pre-trained model, typically creating a specialized static configuration for a specific downstream task. NCNs, instead, focus on *dynamically* modulating the behavior of the existing (potentially fixed) base model parameters during a single forward pass. The adaptation is driven by the immediate context via the NCN signals, allowing for potential shifts in processing strategy within a sequence or across different inputs, rather than relying on a pre-set, task-specific adapter configuration.
- **Hypernetworks:** Hypernetworks generate weights for a primary network based on some input, such as a task embedding [8, 16]. While context-dependent, they typically produce a *static* set of parameters (e.g., for adapters or specific layers) *before* the main processing of an instance begins. NCNs distinguish themselves by generating *dynamic modulation signals*

(like gain factors or temperature scalars) that influence the primary network's computations *during* the forward pass. These signals can potentially change based on evolving context or internal state, offering a more continuous form of adaptation compared to the pre-computation of static weights by hypernetworks.

- **Gating Mechanisms (LSTM/GRU/FiLM):** Internal gating, as seen in LSTMs [11] and GRUs [3], controls information flow based on the current input and local hidden state. FiLM layers [16] compute affine transformation parameters (gain/bias) based on conditioning information. NCNs can be seen as generalizing these ideas under a neuromodulatory framework. Key differences include: (1) NCNs typically derive their modulation signals from a more *global* context (e.g., pooled sequence representation, task ID) computed by a dedicated parallel network, rather than just local state. (2) NCNs propose a broader range of potential modulation targets inspired by biology, including attention temperature and activation function shapes, beyond the gain/bias or flow control typical of these other gating mechanisms.
- **Controllable Generation Techniques:** Methods like CTRL [13] use control codes in the prompt, while PPLM [4] modifies activations during decoding. NCNs offer a different approach by providing an *architectural substrate* for control. Desired attributes (style, safety, topic) could be fed as external inputs to the NCN [Section 3.2], allowing it to generate internal modulation signals that directly shape the model's processing (e.g., adjusting attention focus, layer gains) to align with the control objective in a potentially more integrated fashion.

In essence, NCNs represent a distinct approach that emphasizes dynamic, context-driven modification of the *processing characteristics* of shared LLM components, inspired by the global, state-setting role of biological neuromodulators. This contrasts with methods focused on input modification (prompting), routing to static experts (MoE), adding static task-specific modules (adapters), generating static weights (hypernetworks), or purely local gating (LSTM/GRU). NCNs aim to provide a flexible, built-in mechanism for LLMs to adapt *how* they process information on the fly.

7. Future Research Directions

The Neuromodulatory Control Network (NCN) architecture presented in this paper offers a novel, biologically inspired framework for dynamic LLM processing. While this work establishes the conceptual foundation and provides a reference implementation, numerous avenues for future research remain to fully explore, validate, and refine this approach.

1. **NCN Architectures and Complexity:** The reference implementation utilizes a simple MLP

for the NCN [Section 3.3]. Future work should explore more sophisticated NCN architectures. Recurrent models (LSTMs [11], GRUs [3]) could enable the NCN to integrate contextual information over time within longer sequences or dialogues. Small Transformer-based NCNs might capture more complex relationships within the control input features. Graph Neural Networks (GNNs) could be suitable if the NCN input incorporates structured relational information. Investigating the trade-offs between NCN complexity, computational overhead, and the richness of learned modulation patterns is essential.

2. **NCN Input Sources and Integration:** Exploring richer and more dynamic NCN inputs [Section 3.2] is crucial. This includes systematically evaluating the utility of internal LLM state metrics (e.g., activation statistics, attention entropy), incorporating finer-grained task information, and designing effective methods for integrating external control signals. Research could also focus on learning *which* input sources are most relevant for effective modulation, perhaps using attention mechanisms within the NCN itself if it processes multiple input types.
3. **Systematic Investigation of Modulation Targets:** Our proposal outlines several potential modulation targets (attention temperature/gating, layer gain, FFN modulation) [Section 3.5]. Rigorous empirical studies are needed to determine the relative effectiveness and potential synergies of modulating different components. Ablation studies disabling specific modulations (e.g., gain-only vs. temperature-only vs. both) and comparing different mechanisms (e.g., FFN output scaling vs. internal gate modulation) will be critical for understanding which forms of modulation are most beneficial for specific tasks or behaviors.
4. **Advanced Training Regimes:** While end-to-end training is the simplest approach, its effectiveness, particularly regarding gradient flow to the NCN, needs thorough investigation [Section 5.2]. Comparative studies involving Reinforcement Learning (RL) [23], where the NCN acts as an agent optimizing task rewards or behavioral objectives, are warranted, especially if some modulation effects are non-differentiable or require long-term credit assignment. Designing effective auxiliary losses that explicitly encourage meaningful or diverse modulation patterns could also stabilize and improve NCN training.
5. **Application-Specific Validation:** NCN-augmented LLMs should be rigorously tested on tasks specifically designed to require high adaptability [Section 5.1]. Promising application areas include:
 - **Continual Learning:** Can NCNs modulate plasticity or processing to mitigate catastrophic forgetting?
 - **Multi-Task Learning Benchmarks:** Quantify improvements over baselines like shared parameters or adapters on diverse task suites (e.g., GLUE, SuperGLUE).
 - **Controllable Generation:** Evaluate the precision and range of control achievable by feeding external signals into the NCN for attributes like style, sentiment, or safety.

- **Long-form Dialogue/Generation:** Assess if NCNs can adapt processing based on evolving context over extended interactions.
 - **Interactive Agents:** Explore NCNs in agents that need to switch strategies based on environmental feedback.
6. **Synergy with Other Architectures:** Investigating combinations of NCNs with existing techniques could yield novel capabilities. Could NCNs dynamically control the routing probabilities in MoE models [7, 22]? Could they modulate the parameters or activation levels of Adapters [12, 18] for finer-grained task adaptation? Could NCNs govern synaptic plasticity rules inspired by biology, acting as meta-learning controllers [6]?
 7. **Theoretical Analysis:** Developing a better theoretical understanding of how NCN modulation affects the learning dynamics, representation geometry, and stability of Transformer models is important. How does dynamic gain interact with Layer Normalization [1]? How does modulating attention temperature influence gradient flow and optimization landscapes?
 8. **Evaluation of Modulation Effectiveness:** Moving beyond standard task metrics, future work should develop methodologies to specifically assess *whether* the NCN is learning and applying meaningful, context-dependent modulation [Section 5.1, 5.2]. This could involve analyzing modulation signal distributions across diverse inputs/tasks, performing interventions (fixing or ablating signals), and correlating signal values with intermediate network states or final performance outcomes.
 9. **Efficiency and Conditional Computation:** Realizing the potential for dynamic resource allocation [Section 5.1] requires further investigation. Can NCNs learn to effectively gate entire layers or large components by driving gain signals towards zero in relevant contexts? Developing training methods that explicitly encourage such sparsity while maintaining performance is a key challenge. Measuring actual reductions in FLOPs or latency in such scenarios is needed.
 10. **Interpretability and Visualization:** Understanding the learned behavior of NCNs is crucial for trust and further development [Section 5.2]. Research is needed into techniques for visualizing NCN signal dynamics, correlating signals with interpretable input features or behavioral modes, and potentially methods to encourage more interpretable modulation strategies (though this may trade off with performance).

By addressing these research directions, the community can build upon the initial NCN concept to develop more robust, adaptable, controllable, and potentially more efficient large language models.

8. Conclusion

Large Language Models have demonstrated remarkable capabilities, yet their predominantly static processing mechanisms limit their ability to dynamically adapt to varying contexts, tasks, and control objectives in a manner analogous to biological intelligence. This paper introduced Neuromodulatory Control Networks (NCNs), a novel architectural approach inspired by the neuromodulatory systems of the brain. NCNs are proposed as small, parallel networks that compute context-dependent signals to dynamically modulate the core computational components of a base LLM, influencing *how* information is processed during a forward pass.

We have detailed the conceptual framework of NCNs, outlining potential input sources, network architectures, and specific mechanisms for modulating attention, layer gains, and feed-forward networks. Furthermore, we presented key details of our open-source PyTorch implementation, providing a concrete realization of the architecture, including custom components necessary for features like attention temperature modulation, to facilitate further research and development by the community. By comparing NCNs to existing approaches like MoE, Adapters, and Hypernetworks, we highlighted the unique focus of NCNs on dynamic, broadcast modulation of shared parameters.

While this work primarily presents the architectural proposal and codebase, we anticipate that NCNs hold significant potential for enhancing LLM adaptability, multi-task performance, controllability, and potentially computational efficiency through dynamic resource allocation. However, realizing this potential requires addressing key challenges related to stability, training dynamics, interpretability, and overhead, which form the basis for extensive future research.

In conclusion, Neuromodulatory Control Networks offer a compelling, biologically grounded direction for enhancing LLM flexibility and context-awareness. We believe that exploring principles derived from biological neuromodulation, as embodied in the NCN architecture, is a promising path toward developing more sophisticated, adaptable, and ultimately more capable AI systems. We release our codebase with the hope of stimulating further investigation and empirical validation of this approach within the research community.

References

- [1] Ba, J. L., Kiros, J. R., & Hinton, G. E. (2016). Layer normalization. *arXiv preprint arXiv:1607.06450*.
- [2] Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., ... & Amodei, D. (2020). Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*.
- [3] Cho, K., Van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., & Bengio, Y. (2014). Learning phrase representations using RNN encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*.
- [4] Dathathri, S., Madotto, A., Lan, J., Hung, J., Frank, E., Molino, P., ... & Fung, P. (2019). Plug and play language models: A simple approach to controlled text generation. *arXiv preprint arXiv:1912.02164*.
- [5] Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- [6] Doya, K. (2002). Metalearning and neuromodulation. *Neural networks*, 15(4-6), 495-506.
- [7] Fedus, W., Zoph, B., & Shazeer, N. (2021). Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *arXiv preprint arXiv:2101.03961*.
- [8] Ha, D., Dai, A., & Le, Q. V. (2016). Hypernetworks. *arXiv preprint arXiv:1609.09106*.
- [9] He, J., Zhou, C., Ma, X., Berg-Kirkpatrick, T., & Neubig, G. (2021). Towards a unified view of parameter-efficient transfer learning. *arXiv preprint arXiv:2110.04366*.
- [10] Hendrycks, D., & Gimpel, K. (2016). Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*.
- [11] Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural computation*, 9(8), 1735-1780.
- [12] Houshy, N., Giurgiu, A., Jastrzebski, S., Morrone, B., De Laroussilhe, Q., Gesmundo, A., ... & Gelly, S. (2019). Parameter-efficient transfer learning for NLP. *arXiv preprint arXiv:1902.00751*.
- [13] Keskar, N. S., McCann, B., Varshney, L. R., Xiong, C., & Socher, R. (2019). Ctrl: A conditional transformer language model for controllable generation. *arXiv preprint arXiv:1909.05858*.

- [14] Lepikhin, D., Lee, H., Xu, Y., Chen, D., Firat, O., Huang, Y., ... & Chen, M. (2020). Gshard: Scaling giant models with conditional computation and automatic sharding. *arXiv preprint arXiv:2006.16668*.
- [15] Mahabadi, R. K., Ruder, S., & Henderson, J. (2021). Parameter-efficient multi-task fine-tuning for transformers via shared hypernetworks. *arXiv preprint arXiv:2106.04489*.
- [16] Perez, E., Strub, F., De Vries, H., Dumoulin, V., & Courville, A. (2018). Film: Visual reasoning with a general conditioning layer. *arXiv preprint arXiv:1709.07871*.
- [17] Pfeiffer, J., Vulic, I., Gurevych, I., & Ruder, S. (2020). AdapterHub: A framework for adapting transformers. *arXiv preprint arXiv:2007.07779*.
- [18] Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., & Sutskever, I. (2019). Language models are unsupervised multitask learners. *OpenAI blog*, 1(8), 9.
- [19] Schultz, W., Dayan, P., & Montague, P. R. (1997). A neural substrate of prediction and reward. *Science*, 275(5306), 1593-1599.
- [20] Shazeer, N. (2020). Glu variants improve transformer. *arXiv preprint arXiv:2002.05202*.
- [21] Shazeer, N., Mirhoseini, A., Maziarz, K., Davis, A., Le, Q., Hinton, G., & Dean, J. (2017). Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. *arXiv preprint arXiv:1701.06538*.
- [22] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention Is All You Need. *arXiv preprint arXiv:1706.03762*.
- [23] Ziegler, D. M., Stiennon, N., Wu, J., Brown, T. B., Radford, A., Amodei, D., ... & Irving, G. (2019). Fine-tuning language models from human preferences. *arXiv preprint arXiv:1909.08593*.

Acknowledgments

The author acknowledges the significant contribution of Google's Gemini 2.5 Pro in generating the initial concept for Neuromodulatory Control Networks (NCNs) and drafting the manuscript text. The author directed the generation process, curated the outputs, and takes full responsibility for the final content of this paper.