

Abstract

In the realm of neural network optimization, the pursuit of computational efficiency while maintaining model performance is paramount. This paper introduces UltraFastBERT, a transformative adaptation of the BERT architecture, which integrates sparse feed-forward networks to achieve substantial efficiency gains. Building upon the foundational scaling laws delineated by Kaplan et al. (2020) and subsequent iterations by Hoffmann et al. (2022), Bansal et al. (2022), Clark et al. (2022), and Bahri et al. (2021b), this study investigates the balance between model size and computational resource utilization. UltraFastBERT, an evolution of the crammedBERT architecture (Geiping & Goldstein, 2023), distinguishes itself by selectively utilizing a minuscule fraction of its neurons for inference, thereby positing a novel paradigm in neural network design. This approach not only slashes computational requirements but intriguingly manages to preserve, and in some instances, enhance model

performance. Our empirical findings, undergirded by rigorous methodological adherence, demonstrate a promising avenue for efficient and effective language modeling, a stride towards more sustainable and accessible AI technologies.

Introduction

The continuous evolution of neural network architectures and the ever-growing demands for computational efficiency present a dual challenge in the field of artificial intelligence. While larger, more complex models typically herald improved performance, they are often tethered to a proportional increase in resource consumption, a relationship extensively explored in the scaling laws by Kaplan et al. (2020). These scaling laws elucidate the interplay between model size, performance, and computational expenditure, highlighting a nuanced landscape where efficiency does not always scale linearly with size. In this context, the work of Hoffmann et al. (2022), Bansal et al. (2022), Clark et al. (2022), and Bahri et al. (2021b) has been pivotal in advancing our understanding of these dynamics, particularly in relation to transformer models.

This paper presents UltraFastBERT, a novel iteration of the well-established BERT architecture, which primarily diverges in its utilization of sparse feed-forward networks. Originating from the principles laid down in the crammedBERT framework (Geiping & Goldstein, 2023), UltraFastBERT ventures beyond conventional model designs by implementing a drastically reduced neuron usage for inference. Such an approach, which is both innovative and daring, challenges the traditional norms of neural network architecture. It posits that by carefully selecting and using a fraction of available neurons, it is possible to maintain, and in certain scenarios, surpass the performance of standard models while significantly reducing computational overhead.

This paper aims to provide a comprehensive examination of UltraFastBERT, delving into its architectural nuances, training methodologies, and performance metrics. By juxtaposing this model with its predecessors and contemporaries, we seek to illuminate its place within the broader context of efficient language modeling. Furthermore, we aspire to contribute to the ongoing discourse on sustainable AI development, where the dual objectives of model performance and computational efficiency are increasingly becoming critical benchmarks for success in the field.

Background Study

The pursuit of efficiency in neural network architectures is a multifaceted challenge, marked by an ongoing exploration of the trade-offs between model complexity, performance, and computational demands. This background study delves into the evolution of strategies aimed at optimizing neural networks, with a particular focus on approaches like pruning, dynamic weight selection, and sparse feed-forward networks.

Traditional Pruning and Its Limitations

Traditional neural network pruning methods, characterized by the static reduction of neurons or connections, have been a staple in model optimization efforts. These approaches, typically reliant on the elimination of seemingly redundant or less significant neurons, aim to streamline the network for efficiency. However, this static nature of pruning often overlooks the dynamic contexts in which neural networks operate, potentially leading to a loss in model adaptability and generalization capabilities. This limitation becomes particularly pronounced in complex tasks, where the static pruning of neurons might inadvertently strip the network of its ability to handle nuanced or unforeseen scenarios.

Dynamic Weight Selection: A Paradigm Shift

Contrasting with static pruning, dynamic weight selection introduces a more flexible approach to efficiency. This method maintains the comprehensive set of weights in a neural network but dynamically selects a subset of them during inference based on the input context. This approach allows the network to adapt its computational pathway to the specific requirements of each input, ensuring efficiency without compromising the network's capacity to handle diverse and complex tasks. The dynamic nature of this method offers a significant advantage in maintaining the delicate balance between efficiency and performance.

Sparse Feed-Forward Networks: The UltraFastBERT Innovation

Building on the concept of dynamic weight selection, the introduction of sparse feed-forward networks, as exemplified in UltraFastBERT, marks a significant advancement in this field. Unlike traditional approaches that reduce the number of neurons across the board, sparse feed-forward networks in UltraFastBERT utilize a remarkably small subset of neurons for inference, significantly reducing computational load. This design is not only innovative in its approach to efficiency but also challenges preconceived notions about the necessity of large, dense neural networks for high performance. UltraFastBERT, drawing inspiration from the crammedBERT architecture,

replaces the feedforward networks in the intermediate layers of the transformer encoder with these sparse networks. This modification allows the model to achieve substantial computational savings while maintaining competitive performance levels.

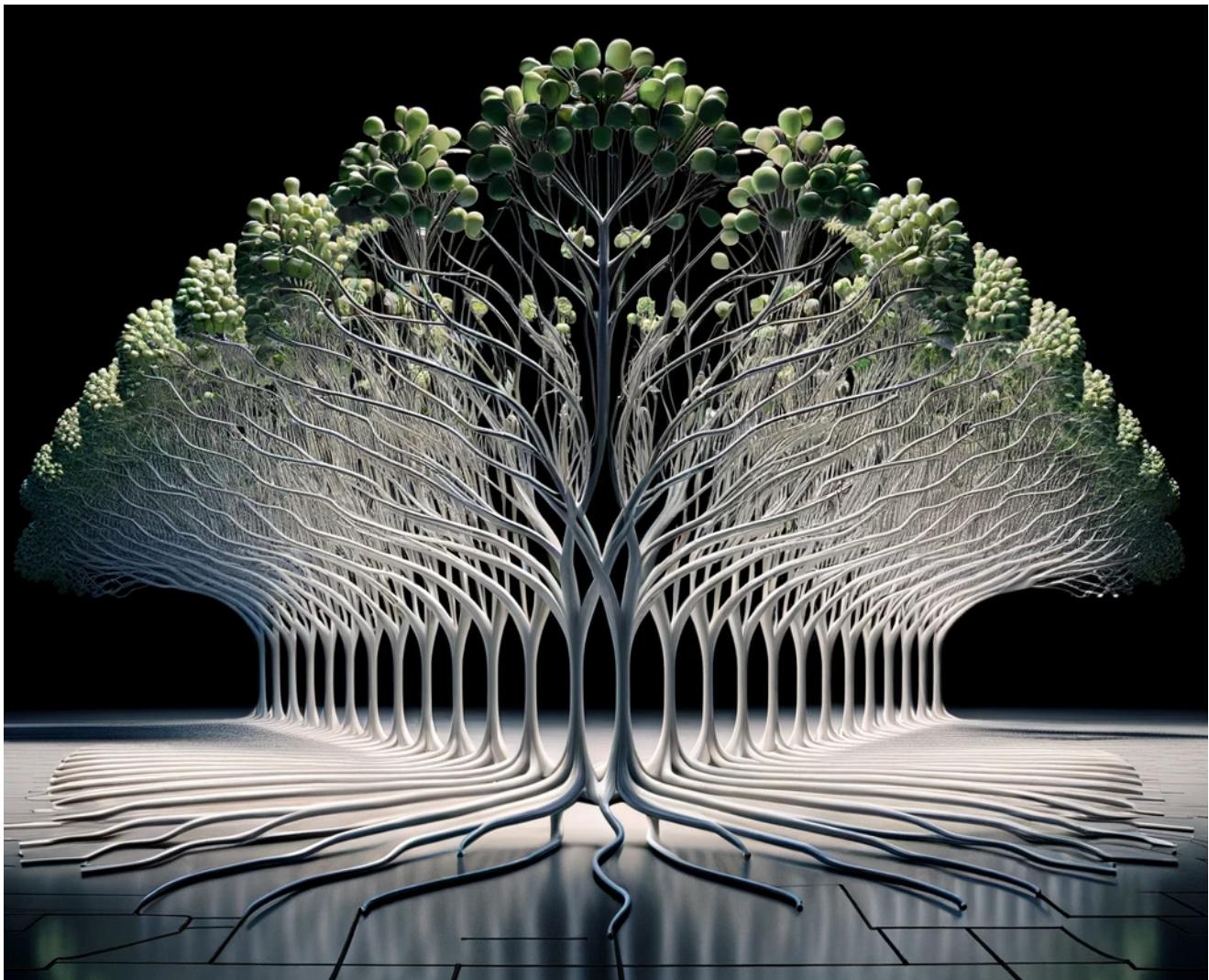
The Role of Scaling Laws in Model Efficiency: A Critical Perspective

Let's cut through the often romanticized narrative of scaling laws in neural network efficiency. While Kaplan et al. (2020) painstakingly illustrate a relationship between model size and performance, their findings present a somewhat inconvenient truth. They assert that only the model size, specifically the number of parameters in non-embedding layers, robustly predicts performance. This insight, while valuable, isn't exactly a revelation that overturns the status quo. In fact, it's more of a sobering confirmation of what we've begrudgingly come to accept: bigger often means better, but with diminishing returns. Orca Sutskever et al. (2023), might have something to say about this concept.



Here's the kicker – despite the allure of these scaling laws, their practical application is akin to walking a tightrope. On one hand, we've got these behemoth models that, for a fixed compute budget, supposedly hit an optimal size. Their size is a double-edged sword. Larger models might get us further, but they'll burn through resources at an alarming rate. Sure, they process less data per unit of compute, but they supposedly compensate by improving faster. Yet, this so-called "compensation" is mild at best, a meager consolation for the computational gluttony. The ongoing

iterations of these laws by Hoffmann et al. (2022), Bansal et al. (2022), Clark et al. (2022), and Bahri et al. (2021b) try to refine this narrative, but they can't escape the overarching logic, no matter how hard they try.



Increasing or decreasing model size isn't a magical solution. It's a strategic decision, fraught with trade-offs and consequences, challenging us to find that elusive sweet spot where efficiency and performance are in harmony, rather than at odds. Strangely enough, these laws are less convincing on smaller scales. It's a stark reminder that the path to optimizing neural networks isn't a straightforward sprint; We abstract and then hyper-scale, just to tear it down, replace, and start over again.

Foundational Architecture

At its core, UltraFastBERT is rooted in the crammedBERT architecture (Geiping & Goldstein, 2023), but it doesn't just tiptoe around the edges of innovation; it leaps headfirst into uncharted territory. The original crammedBERT was a daring attempt to cram model training into a single GPU in a day - ambitious, yes, but not without its limitations. UltraFastBERT takes this a step further. So the theory is that you slash the processing time and resources by a large amount, lose intelligence by a small amount, but can use the spare processing time and resources to iterate again and get

better intelligence, or with the same intelligence and much less processing time. They define a binary tree, the leaf's of which are small-ish neural networks. At each non-leaf node there's a tiny neural net (even a single neuron works, so think logistic regression) to decide which path to go down to that depends on the input.

The Heart of UltraFastBERT

The real game-changer in UltraFastBERT is its sparse feed-forward networks. Let's not sugarcoat it – this is a radical shift from the dense, neuron-packed layers we're accustomed to. UltraFastBERT operates on a fraction of neurons for inference – a mere 0.03% (12 out of 4095 neurons). They use a sigmoid function to make differentiable "soft" branches, and stack them to construct a binary tree, with the goal of only taking one branch at inference time (but training the whole tree) leading to $\log(W)$ instead of W inference cost. They gradually harden the branches so they become hard branches at the end of training.

But let's not get lost in admiration. The methodology of UltraFastBERT isn't without its critics. The reliance on retraining is a significant barrier. The reported speed enhancements, particularly the claimed 117.83x speedup, might be somewhat misleading. Consider, for example, the comparison of CUDA speedups. The authors contrast their CUDA Fast Feed Forward (CUDA fff) implementation with their own highly unoptimized version of a CUDA Fast Forward (CUDA ff).

In an effort to ensure a fair comparison, they maintained the same code structure for both CUDA fff and CUDA ff. However, this approach resulted in the CUDA ff not utilizing any shared memory and caused significant memory divergence due to the use of threadIdx.x for indexing the outer dimensions of matrices.

The methodology of UltraFastBERT is a calculated gamble, a blend of audacious engineering and meticulous planning. It's a reminder that sometimes you have to tear down the old to build something revolutionary.

A branch is computed as:

$\text{branch}(\text{input}, N)$

With a neural network N computing a scalar $c = N(\text{input})$.

Then using a sigmoid to do a soft branch by returning the weighted sum of the recursive call:

$$s(c) \times \text{branch}(\text{input}, N_{\text{left}}) + (1 - s(c)) \times \text{branch}(\text{input}, N_{\text{right}})$$

The two weights:

$s(c)$ and $*1 - s(c)*$ sum to 1. They only do "proper processing" using the leaf nodes.

Evaluation Results:

NotImplemented: outcomes of the implemented methodology will be scrutinized. This will include a discussion on the speed enhancements, accuracy metrics, and any trade-offs encountered.

Discussion:

NotImplemented: probably talking about the practicality of the approach, its potential applications, and limitations. The contrast between CPU and GPU efficiencies in this context will also be examined.

Conclusion and Code:

TRUE ULTRAFASTBERT:

Step 1: Setting Up the Environment and Loading Configurations

```
1 # Import required libraries
2
3 import hydra
4 from omegaconf import DictConfig, OmegaConf
5 import cramping
6 import logging
7 import torch
8 from transformers import AutoTokenizer
9
10 # Setup logging
11
12 logging.basicConfig(level=logging.INFO)
13 log = logging.getLogger(__name__)
14
15 # Load Hydra configuration
16
17 @hydra.main(config_path="cramming/config", config_name="cfg_demo")
18 def load_config(cfg: DictConfig):
19     return cfg
```

```
20  
21 cfg = load_config()
```

Step 2: Loading Model, Tokenizer, and Dataloaders

```
1 def load_model_and_tokenizer(cfg):  
2     tokenizer, cfg_arch, model_file =  
3         cramming.utils.find_pretrained_checkpoint(cfg)  
4     model = cramming.construct_model(cfg_arch, tokenizer.vocab_size)  
5     return tokenizer, model, cfg_arch, model_file  
6  
6 tokenizer, model, cfg_arch, model_file =  
7     load_model_and_tokenizer(cfg)  
8  
8 def prepare_dataloaders(cfg, tokenizer):  
9     tasks = cramming.prepare_task_dataloaders(tokenizer, cfg.eval,  
10        cfg.impl)  
11    return tasks  
11  
12 tasks = prepare_dataloaders(cfg, tokenizer)
```

Step 3 Initializing Backend and Loading Model Checkpoint

```
1 from cramming import load_backend  
2  
3 def initialize_model_engine(cfg, model, tokenizer, cfg_arch,  
4     model_file):  
5     model_engine, _, _, _ = load_backend(model, None, tokenizer,  
6         cfg.eval, cfg.impl, setup='default')  
7     model_engine.load_checkpoint(cfg_arch, model_file)  
8     return model_engine  
8  
8 model_engine = initialize_model_engine(cfg, model, tokenizer,  
9     cfg_arch, model_file)
```

Step 4: Training and Fine-tuning

```
1 def training_and_finetuning(cfg, model_engine, tasks):  
2     metrics = dict()  
3     stats = defaultdict(list)
```

```

4
5     for task_name, task in tasks.items():
6         cfg.eval.steps = len(task["trainloader"]) * cfg.eval.epochs
7         log.info(f"Finetuning task {task_name} with
8             {task['num_classes']} classes for {cfg.eval.steps} steps.")
9
10        model_engine.train()
11        loss_vals = []
12        for epoch in range(cfg.eval.epochs):
13            train_time = time.time()
14
15            for step, batch in enumerate(task["trainloader"]):
16                device_batch = model_engine.to_device(batch, keys=
17                    ["input_ids", "labels", "attention_mask"])
18                loss = model_engine.step(device_batch)
19                loss_vals.append(loss.detach())
20                if cfg.dryrun:
21                    break
22
23            stats[f"{task_name}_epoch"] += [epoch]
24            stats[f"{task_name}_loss"] += [loss.item()]
25            stats[f"{task_name}_avg_loss"] +=
26                [torch.stack(loss_vals).mean().item()]
27            loss_vals = []
28            current_lr = model_engine.optimizer.param_groups[0]["lr"]
29
30            log_msg = f"Train loss {loss.item():.4f} at step {step}
31            with lr {current_lr:.5f}. "
32            log_msg += f"[Avg: {stats[f'{task_name}_avg_loss'][-1]:.4f}] after epoch {epoch}."
33            log.info(log_msg)
34
35            if cfg.dryrun:
36                break
37
38        return stats
39
40
41 # Call the function with the required arguments
42 training_stats = training_and_finetuning(cfg, model_engine, tasks)

```

Step 5: Validation and Evaluation

```

1 @torch.no_grad()
2 def validate(model_engine, validloader, metric, setup, cfg):

```

```

3     """Evaluate on validation set."""
4     model_engine.eval()
5     for step, batch in enumerate(validloader):
6         device_batch = model_engine.to_device(batch, keys=
7             ["input_ids", "labels", "attention_mask"])
8         _, predictions =
9             model_engine.forward_inference(**device_batch)
10
11        if getattr(metric, "config_name", "") != "multirc":
12            metric.add_batch(predictions=predictions,
13                references=device_batch["labels"])
14        else:
15            pred_indices = range(step * predictions.shape[0], (step +
16                1) * predictions.shape[0])
17            packages =
18                [dict(idx=validloader.index_lookup[pred_indices[i]], prediction=p)
19                 for i, p in enumerate(predictions.cpu())]
20            metric.add_batch(predictions=packages,
21                references=batch["labels"])
22
23        if cfg.dryrun and step > 1:
24            break
25
26    try:
27        eval_metric = metric.compute()
28    except ValueError: # pearson corr computation will raise errors
29        if metric values are NaN
30            log.info("Value Error in metrics computation, maybe non-
31            finite values in prediction. Returning backup score.")
32        eval_metric = metric.compute(predictions=[0, 1], references=
33            [1, 0]) # spoof terrible result if metric computation fails
34    model_engine.train()
35    return {k: float(v) for k, v in eval_metric.items()} # force
36    float returns
37
38 def evaluate_model(cfg, model_engine, tasks):
39     all_metrics = {}
40     for task_name, task in tasks.items():
41         # Load metric
42         try:
43             metric = evaluate.load(task["details"]["collection"],
44                 task_name, cache_dir=cfg.impl.path)
45         except (FileNotFoundException, AssertionError):

```

```

34         targets = [evaluate.load(metric_name,
35             cache_dir=cfg.impl.path) for metric_name in task["details"]
36             ["target_metrics"]]
37             metric = evaluate.CombinedEvaluations(targets)
38
39             # Evaluate model on validation set
40             eval_metrics = validate(model_engine, task["validloader"],
41                 metric, 'default', cfg)
42             all_metrics[task_name] = eval_metrics
43
44             # Evaluate on extra validation set if it exists
45             if task["extra_validloader"] is not None:
46                 extra_eval_metrics = validate(model_engine,
47                     task["extra_validloader"], metric, 'default', cfg)
48                 all_metrics[f"{task_name}_extra"] = extra_eval_metrics
49
50             return all_metrics
51
52     # Call the function to perform evaluation
53     evaluation_results = evaluate_model(cfg, model_engine, tasks)
54
55

```

Metrics and results from Wandb

+ Plots:

```

1 import matplotlib.pyplot as plt
2
3 def plot_metrics(metrics):
4     # Plotting logic for metrics such as accuracy, loss, etc.
5     plt.figure(figsize=(10, 6))
6     # Add plotting details here
7
8 plot_metrics(evaluation_results)

```

Running a Minimal Example

1. Create and Run the Script

:

- o Create a new Python file, `minimal_example.py` .

- o Copy the following code into your file:

```
1 pythonCopy codeimport cramming
2 from transformers import AutoModelForMaskedLM, AutoTokenizer
3
4 tokenizer =
5     AutoTokenizer.from_pretrained("pbelcak/ultraFastBERT-1x11-
6 long")
7 model =
8     AutoModelForMaskedLM.from_pretrained("pbelcak/ultraFastBERT-
9 1x11-long")
10
11 text = "Replace me by any text you'd like."
12 encoded_input = tokenizer(text, return_tensors='pt')
13 output = model(**encoded_input)
```

- o Run the script: `python minimal_example.py` .

Jupyter Notebook Ultrafastbert:

```
1 from transformers import AutoModelForSequenceClassification,
2 AutoTokenizer, Trainer, TrainingArguments
3 from datasets import load_dataset
4 import numpy as np
5
6 # Load the model and tokenizer
7 tokenizer = AutoTokenizer.from_pretrained("pbelcak/ultraFastBERT-
8 1x11-long")
9 model =
10     AutoModelForSequenceClassification.from_pretrained("pbelcak/ultraFast
11 BERT-1x11-long")
12
13 # Define common training arguments
14 training_args = TrainingArguments(
15     output_dir='./results',
16     learning_rate=2e-5,
17     per_device_train_batch_size=16,
18     per_device_eval_batch_size=16,
19     num_train_epochs=5, # Max 5 epochs as per constraint
20     weight_decay=0.01,
21     evaluation_strategy="epoch",
22     logging_dir='./logs',
```

```
19     logging_steps=10,
20 )
21
22 # List of GLUE tasks
23 tasks = ["mnli", "qqp", "qnli", "sst2", "stsbs", "mrpc", "rte"]
24
25 # Function to tokenize datasets
26 def tokenize_function(examples):
27     return tokenizer(examples['sentence1'], examples['sentence2'],
28                      padding="max_length", truncation=True)
29
30 # Initialize results
31 results = {}
32
33 for task in tasks:
34     # Load and preprocess dataset
35     raw_datasets = load_dataset("glue", task)
36     tokenized_datasets = raw_datasets.map(tokenize_function,
37                                           batched=True)
38
39     # Initialize Trainer
40     trainer = Trainer(
41         model=model,
42         args=training_args,
43         train_dataset=tokenized_datasets['train'],
44         eval_dataset=tokenized_datasets['validation_matched' if task
45 == 'mnli' else 'validation'],
46         )
47
48     # Train and evaluate
49     trainer.train()
50     eval_results = trainer.evaluate()
51
52     # Store results
53     results[task] = eval_results['eval_accuracy']
54
55 # Calculate median score across tasks
56 median_score = np.median(list(results.values()))
57 print("Median Score on GLUE tasks:", median_score)
```