

# ResNet-BK コードベース暫定レビュー

## 1. コードレビュー

### 1-1. 作業内容の把握 (What / 何をするコードか)

- リポジトリ全体の役割：このリポジトリは  $O(N)$  言語モデルアーキテクチャ (ResNet-BK) を開発し、Transformer に依存しない大規模言語モデルを作る試みである。 README によると、BK コアの導入によりテキスト長に対する計算量を線形化し、スパース化と量子化を組み合わせてトレーニングコストの数百万倍削減を目指す <sup>1</sup>。
- 主要ディレクトリとファイル
  - `1_BK_Language_Model_PoC/` - 言語モデル実装と実験結果。
    - `src/models/resnet_bk.py` - ResNetBKBlock と LanguageModel の実装。トークン埋め込みと位置埋め込みを持ち、BKCoresFunction による  $O(N)$  resolvent を通して出力特徴を得る。MoE 層を挟んで最終的に線形ヘッドへ投影する <sup>2</sup>。
    - `src/models/bk_core.py` - BK コア。三対角行列の逆行列の対角要素を求めるために順方向と逆方向の漸化式を使う。解析的な勾配と数値的な勾配をブレンドする GRAD\_BLEND パラメータを持ち、数値安定化のために勾配をクリップしている <sup>3</sup>。
    - `src/models/moe.py` - トップ-k ルーティングを用いた Sparse Mixture-of-Experts レイヤ。ゲーティングネットワークにより K 個の専門家出力を選び、スパース化を実現する <sup>4</sup>。
    - `src/models/quantized_bk_core.py` - 量子化対応 BK コア。校正により FP32 の重みを 8bit 整数にマッピングし、学習時にはフェイク量子化を行い推論時に INT8 に変換する <sup>5</sup>。
    - `src/models/configurable_resnet_bk.py` - 多数のフラグを持つ ResNetBKConfig dataclass。解析勾配(Phase2)、モエ層、量子化(Phase4)、勾配チェックポイント、カスタム CUDA カーネル(Phase5)、Adaptive Computation Time などが有効/無効にできる <sup>6</sup>。
    - `train.py` - 言語モデルのトレーニングスクリプト。ConfigurableResNetBK をインスタンス化し、オプティマイザやスケジューラを設定してメトリクスをロギングする <sup>7</sup>。
  - `2_Scaling_Benchmarks/` - スケーリング実験やベンチマーク用スクリプト。`run_scaling_experiments.py` はさまざまなモデルサイズをプロファイルし、処理速度・パラメータ数のスケーリング法則を探る <sup>8</sup>。 `run_pile_benchmark.py` は Pile データセットでの比較実験を行う <sup>9</sup>。

### 1-2. 進捗の把握 (Where / どこまで来ているか)

- Step1: アーキテクチャ刷新 - 完了。BK コアをベースとした ResNet-BK モデルと MoE 層が実装されており、Transformer baseline と比べて同容量で  $6.7 \times$  高速 ( $N=2048$ ) かつ学習能力を示した <sup>1</sup>。
- Step2: アルゴリズム革新 (解析勾配) - 部分実装。`bk_core.py` では解析的な勾配と仮説 <sup>7</sup>に基づくハイブリッド勾配を混合しており、grad\_blend パラメータのグリッドサーチを行った結果、Perplexity が 512 シーケンス長で 428 と Transformer baseline より大幅に改善した <sup>10</sup>。
- Step3: スパース化 (MoE) - 実装済み。`moe.py` で top-k ルーティングの Mixture-of-Experts を提供し、モデル計算をスパース化している <sup>4</sup>。
- Step4: 圧縮 (量子化・剪定・蒸留) - 試験的実装。量子化対応 BK コア、複素数量子化、INT4 MoE、構造化剪定、知識蒸留を組み合わせたパイプラインで約  $100 \times$  圧縮を達成したと README が報告している <sup>11</sup>。
- Step5: ハードウェア最適化 (カスタム CUDA/Triton、勾配チェックポイント) - `ResNetBKConfig` にフラグはあるが未実装。README では今後の課題として GPU カーネルの最適化やメモリ節約技法を挙げている。

- Step6: アルゴリズムによる計算削減 (Adaptive Computation Time、スパース化強化) – フラグのみ存在。早期退出や多重スケール処理、学習されたスパース化などが検討課題。
- Step7: 学習理論革新 (Koopman operator) – 将来計画として、バックプロパゲーションを Koopman 作用素に基づく物理学的学習に置き換えることが提案されている<sup>12</sup>。
- ベンチマーク結果 – README によれば、WikiText-2 での小規模ベンチマークで ResNet-BK は Transformer baseline より低 perplexity (590→428) を達成し、T4 GPU 1枚でも 512 長のシーケンスで OOM せず訓練できる<sup>13</sup>。<sup>1</sup> では 6.7×速いと報告されているが、大規模データセットでの実験は未実施。

### 1-3. 内容の詳細 (How / どう書かれているか)

- 数学コアとアルゴリズム – `bk_core.py` の BK コアは三対角行列の逆行列の対角要素を求める  $O(N)$  アルゴリズムで、順方向 ( $\theta$ ) と逆方向 ( $\phi$ ) の漸化式を計算し、それらを組み合わせて Resolvent の対角成分を得る<sup>3</sup>。解析的な勾配を導出しつつ、数値的安定性のため勾配クリップし、勾配ブレンド係数 `GRAD_BLEND` を導入して理論勾配と実験的勾配を混合する点が特徴である。
- トレーニング設計 – `train.py` ではモデル定義を `ConfigurableResNetBK` に委譲し、AdamW オプティマイザや線形ウォームアップ付きコスサイン減衰スケジューラを使用している<sup>7</sup>。学習ループはデータセットからバッチを取り出し、損失・perplexity・速度を記録しながらエポックを進める。`ConfigurableResNetBK` は多数のフラグに基づいて自動的に量子化、剪定、蒸留、勾配チェックポイント等を適用する。
- 実装上の工夫・懸念
- スパース MoE : `moe.py` では Gumbel-Softmax または top-k 関数で専門家を選択し、選択された専門家の出力を正規化して合成する。専門家ごとに重みを学習しつつ、ルーティングの負荷を減らす<sup>4</sup>。
- 量子化 : `quantized_bk_core.py` では重みを偽量子化する関数 `fake_quantize_per_tensor` と `fake_quantize_per_channel` を使い、校正後に INT8 表現に変換して推論時のメモリを大幅に節約する<sup>5</sup>。ただし Conv/Linear など一般的なレイヤの量子化は未実装で、MoE と BK コアに限定されている。
- `ConfigurableResNetBK` : `ResNetBKConfig` はフラグが多岐に渡るため、設定ミスが起きやすい。Preset 構成が `get_default_config` などにまとめられているが、コメントが少なく意図が読み取りにくい。<sup>6</sup> に記述されたフラグは、GradBlend、Koopman 勾配、量子化、剪定、カスタムカーネル、勾配チェックポイント、ACT など計 30 個近くある。
- 数値安定性 : BK コアは勾配計算が不安定になりやすいため、クリップ幅や正則化パラメータのチューニングが重要。`GRAD_BLEND` の最適値を決めるためのグリッドサーチが行われているが、数学的根拠の説明が不足している。
- エラー・警告 : `quantized_bk_core.py` の校正関数では、range 推定のためにデータセット全体を一度通す必要があり、メモリ消費が大きい。また、`train.py` は DistributedDataParallel を想定していないため、大規模データセットでの拡張性が課題。

### 1-4. 評価 (So what / どれくらい価値があるか)

- 新規性 (Novelty) – 三対角行列の逆行列の対角要素を利用して全結合でない言語モデルを実現する BK コアは斬新で、解析的勾配とのハイブリッド化や MoE によるスパース化など新しい試みが多い。<sup>3</sup> に記載されるアルゴリズムは高速であり、既存 Transformer と比べて理論的・実験的に優位性を示している。
- 工学的完成度 (Engineering quality) – コードは POC 段階のため、テストやコメントが少なく、設定項目が多い割に利用方法が不透明である。一方で `ConfigurableResNetBK` によって多数の技術要素を切り替えられる拡張性を確保している。ベンチマークスクリプトやデータローダが整備されており、Colab 上で結果を再現しやすい。
- 実験の厳密さ (Experiment rigor) – README では perplexity や速度のベンチマーク結果が示されているが、データセット規模が小さい (WikiText-2 のみ) ため、一般化性能を評価するには不十分。解析勾配の効果や量子化の影響に関する詳細なアブレーションが欠如している。

- ・**将来の拡張性 (Extensibility)** - `ResNetBKConfig` により Step5～Step6 の機能を簡単に追加できる設計が魅力である。量子化/剪定/蒸留パイプラインをさらに他のモジュールへ適用する余地がある。  
Koopman 作用素ベースの学習や物理モデルとの統合など研究余地が大きい。<sup>12</sup>

**大学 1 年生としての相対評価**：大学 1 年生がこのコードベースを扱うには難易度が高いが、読み進めればシステム全体の流れを理解しやすい構造になっている。学習アルゴリズムや数値計算の基礎を履修済みなら、BK コアや量子化の仕組みを学ぶ良い教材になる。

**プロ視点で見た成熟度**：POC 段階の研究プロジェクトであり、新しいアイディアを試すための土台としては優秀だが、プロダクションレベルの信頼性や大規模実験の検証は不足している。性能を本格的に検証するにはデータセット拡大やより詳しい ablation が必要である。

※ このレビューは現在閲覧しているコードと README に基づく暫定的なものであり、今後のコミットやドキュメント更新により内容が変わる可能性がある。

## 2. リサーチタスクマップ

以下では、Step5～Step6 などのリサーチ方向性に基づき、約40項目のタスクを列挙する。ID は `(テーマ記号)-(番号)` で示し、Difficulty は T4 無料枠での難易度を想定する。

### 2-1. リサーチ方向性 1：ハードウェアの限界を超える (Step5)

#### テーマ (A) - カスタム CUDA / Triton カーネル

ID	タイトル	種別	難易度	説明	依存
A-01	T4 向け BK コア単一 CUDA カーネルの PoC	[実装]/[実験]	★★★★☆	BK コアでボトルネックとなる漸化式計算を特定し、CUDA カーネルに書き換えてレジスタと共有メモリを活用する。GPU 最適化論文によると、計算依存を減らし共有メモリで再利用することで 6×高速化が期待できる <sup>14</sup> 。	None
A-02	Triton による ResNet-BK ブロックの自動カーネル生成	[実装]	★★★★☆	Triton 言語で BK ブロックと MoE を記述し、auto-tuner で T4 に最適なスレッドブロックを探索する。Emergent Mind や Stanford の研究では LLM によるカーネル生成が標準関数より高速になることが報告されている。	A-01
A-03	アクティベーション圧縮と shared memory の両立	[理論]/[実装]	★★★☆☆	CUDA 最適化の論文 <sup>15</sup> によれば、計算集約型アルゴリズムではレジスタを節約して共有メモリに値を格納することでスループットを高め、メモリ帯域のボトルネックではデータを再構成してマルチパス処理で高速化できる。BK コアに応用し、shared memory と local memory のトレードオフを検討する。	A-01

ID	タイトル	種別	難易度	説明	依存
A-04	Tensor Cores 活用の検討	[理論]/[実験]	★★★★☆	Mixed-precision 行列演算を Tensor Core で高速化するために、BK コアの計算を半精度/単精度に分離する実験。適切なスケーリングと量子化の組合せで精度を保ちつつ 2-3×高速化が期待される。	A-01

#### テーマ (B) - 勾配チェックポイントとメモリ最適化

ID	タイトル	種別	難易度	説明	依存
B-01	<code>torch.utils.checkpoint</code> による BK コアのメモリ削減	[実装]/[実験]	★★☆☆☆	PyTorch ブログではアクティベーションチェックポイント (AC) が保存テンソル数を減らし、計算とメモリのトレードオフを実現する <sup>16</sup> 。BK コア各層で AC を有効にし、GPU メモリ消費を測定する。	None
B-02	Min-cut パーティショナによる自動再計算領域の探索	[実装]/[実験]	★★★★☆	<code>torch.compile</code> の min-cut partitioner は計算グラフを分割して保存テンソル数を最小化する <sup>17</sup> 。BK モデルを <code>torch.compile</code> で最適化し、パーティショナ結果を可視化する。	B-01
B-03	複数チェックポイント戦略の最適化	[実験]	★★★★☆	AC のメモリ削減効果はチェックポイント数と配置に依存する <sup>18</sup> 。ミニバッチサイズを拡大するため、チェックポイントを増やしたときの速度・メモリの推移を調査する。	B-01
B-04	勾配キャッシング+勾配のメモリ圧縮	[実装]/[理論]	★★★★☆	ACT/AdaGrad のようなアルゴリズムで勾配計算をキャッシュし、過去勾配を圧縮して保存する。BK コアでは過去の $\phi, \theta$ シーケンスを再利用できる可能性がある。PyTorch の gradient accumulation や Windows streaming を参考にする。	None

## 2-2. リサーチ方向性 2：アルゴリズムによる計算削減 (Step6)

### テーマ (C) - Adaptive Computation Time / Early Exit

ID	タイトル	種別	難易度	説明	依存
C-01	Early-Exit ResNet-BK ブロックの設計	[実装]/[実験]	★★★★☆☆	例ごとに必要な層数を動的に決定する。PMLR の論文では、入力によって早期に退出できる深層ネットワークにより 2.8× の推論速度向上を達成している <sup>19</sup> 。同様に、ResNet-BK の各ブロック後に classifier を付け confidence が高ければ exit する。	A-02
C-02	ACT 損失の導入とスカラー $\tau$ 学習	[理論]/[実装]	★★★★★☆	Adaptive Computation Time アルゴリズムでは、追加損失を通じて計算ステップ数 (halting score) を学習する。BK モデルに ACT を組み込み、平均計算ステップを制約しながら性能を測定する。	C-01
C-03	早期退出ポリシの強化学習	[理論]/[実験]	★★★★★☆	EXIT ポリシーを layer-wise weighted binary classification として学習する手法 <sup>20</sup> を実装し、推論コストと精度のトレードオフを最適化する。	C-01
C-04	忍耐値 (patience) ベースの早期終了法の比較	[実験]	★★★★☆☆	BERT Loses Patience などの先行研究を参考に、confidence が一定回数続けて高い場合に終了する手法を実装して比較する。	C-01

### テーマ (D) - スパース化の深化 (MoE と bk\_core)

ID	タイトル	種別	難易度	説明	依存
D-01	MoE の Router 温度 制御と負荷バランシング	[実装]/[実験]	★★★★☆☆☆	top-k ルーティングの温度パラメータを適応的に調整し、各専門家の使用率を均等化する。負荷バランス損失を追加し、スパース度と精度の関係を調査する。	C-01
D-02	BK コアへの 学習スパース 化の適用	[理論]/[実装]	★★★★★☆	ネットワーク剪定の文献によると、構造化剪定はチャンネルや層を丸ごと削除することで高速な推論を実現する <sup>21</sup> 。BK コアのパラメータをグループ化して重要度に基づき削減し、精度への影響を評価する。	C-01
D-03	N:M スパース 化規則の導入	[実装]/[実験]	★★★★★☆	GPU に最適化された N:M スパース性（例：2:4）を BK モデルに適用する。CUDA カーネルを対応させることで実際の速度向上を検証する。	D-02

ID	タイトル	種別	難易度	説明	依存
D-04	確率的 MoE とルーティングの自己蒸留	[理論]/[実験]	★★★★☆	MoE ルーターを学習可能な確率的モデルとして扱い、複数のルーティングをサンプリングして平均化する。教師モデルによる自己蒸留を行い、ルーティングの安定性を向上させる。	D-01

#### テーマ (E) - マルチスケール & その他 Step6 技術

ID	タイトル	種別	難易度	説明	依存
E-01	マルチスケール周波数埋め込みの導入	[理論]/[実装]	★★★★☆	ResNet-BK の位置埋め込みを周波数帯域別に複数解像度で表現し、長距離依存性を効率的に捉える。既存の Random Fourier Features などと比較する。	C-01
E-02	マルチスケール BK コア階層の構築	[理論]/[実装]	★★★★☆	小さいレゾルベントを階層的に組み合わせ、長いシーケンスを効率的に処理するアルゴリズムを考案する。各スケール間で再帰的に連携し、メモリと計算を削減する。	E-01
E-03	学習可能な潜在長さの動的縮退	[理論]/[実験]	★★★★☆	長いシーケンスを一部サマリーに圧縮し、BK モデルに入力する Adaptive Token Merging を実装する。モジュールが重要でないトークンを融合し計算を節約する。	C-01

#### 2-3. リサーチ方向性 3：学習理論の革新 (Step2)

##### テーマ (F) - Koopman 作用素とハイブリッド勾配

ID	タイトル	種別	難易度	説明	依存
F-01	Koopman 作用素の基礎理論調査	[ドキュメント]/[理論]	★★☆☆☆	SIAM の解説によれば、Koopman 作用素は非線形力学を線形作用素に移すことで、観測空間上で状態遷移を記述する <sup>22</sup> 。まず有限次元表現や固有関数の性質を学習し、ResNet-BK へ応用するための理論背景を整える。	None
F-02	Koopman 固有関数を学習するニューラルアップローチ	[実装]/[理論]	★★★★☆	文献では固有関数の基底をニューラルネットで近似し、線形作用素を推定する方法が提案されている <sup>23</sup> 。BK モデルの重みを固定し、出力を固有関数空間に投影して学習する実験を行う。	F-01
F-03	ハイブリッド勾配 <code>GRAD_BLEND</code> の理論解析	[理論]/[実験]	★★★★☆	<code>bk_core.py</code> の勾配ブレンドは経験的に選ばれており、理論的な根拠が薄い。安定性と収束性を解析し、Optimal Control 理論や学習率スケジューリングと関連付ける。	F-01

ID	タイトル	種別	難易度	説明	依存
F-04	バックプロパゲーションを不要にする物理学的学習	[理論]	★★★★★	Koopman 作用素やパーシステントホモロジーなどをを利用して勾配無しで学習する手法を探索する。非線形モード分解やスパース回帰を組み合わせ、ResNet-BK を物理学的に訓練できるか検証する。	F-01

#### テーマ (G) - データセット設計・評価指標・蒸留・量子化

ID	タイトル	種別	難易度	説明	依存
G-01	大規模データセットへの拡張とペルプレキシティ評価	[実験]/[ドキュメント]	★★★★☆☆	Pile や Common Crawl などに拡張し、BK モデルと Transformer の perplexity を比較する。ベンチマークの再現性と fairness を確保する。	None
G-02	知識蒸留フレームワークの実装	[実装]/[実験]	★★★★★☆	Lightly 社の解説によると、大きな教師モデルの出力分布を小さな学生モデルに模倣させることで効率的なモデルが得られる <sup>24</sup> 。ResNet-BK を学生として Transformer を教師とする蒸留実験を行う。	G-01
G-03	蒸留損失と圧縮手法の組合せ効果	[実験]	★★★★★☆	蒸留と量子化・剪定を組み合わせた場合の精度と速度を検証する。量子化によりモデルサイズを 4×縮小でき、4bit 量子化では 8×縮小となるが注意が必要 <sup>25</sup> 。	G-02
G-04	構造化剪定アルゴリズムの比較	[実験]	★★★★★☆	Medium の解説では、構造化剪定はフィルタやチャネル全体を削除し、標準ハードウェア上で実際の速度向上を実現する <sup>21</sup> 。さまざまな剪定率で BK モデルを再訓練し、精度と推論速度を測定する。	G-01
G-05	量子化スキームと校正戦略の評価	[実験]	★★★★★☆	Clarifai の記事では、32bit → 8bit 量子化でモデルサイズが 4×縮小し、16×の性能/ワット向上が得られる <sup>25</sup> 。INT4 では 8×縮小だが校正が難しい <sup>26</sup> 。BK モデルで per-tensor/per-channel 量子化を比較し、推論速度と精度を評価する。	G-01
G-06	多言語・長文への拡張	[実験]	★★★★☆☆	日本語や多言語コーパスで BK モデルの性能を評価し、長文 (N>4096) での処理能力を検証する。	G-01

ID	タイトル	種別	難易度	説明	依存
G-07	評価指標の追加 (energy efficiency, latency)	[ドキュメント]/ [実験]	★★★☆☆	電力消費や latency を測定し、量子化・剪定・カスタムカーネル適用後の効率を定量化する。	G-01

## 機械処理向け JSON サマリ（上位 20 件の例）

```
[
  {"id": "A-01", "title": "T4 向け BK コア単一 CUDA カーネルの PoC", "type": ["実装", "実験"], "difficulty": "★★★★☆☆", "dependency": null},
  {"id": "A-02", "title": "Triton による ResNet-BK ブロックの自動カーネル生成", "type": ["実装"], "difficulty": "★★★★☆☆", "dependency": "A-01"},
  {"id": "B-01", "title": "torch.utils.checkpoint による BK コアのメモリ削減", "type": ["実装", "実験"], "difficulty": "★★☆☆☆☆", "dependency": null},
  {"id": "C-01", "title": "Early-Exit ResNet-BK ブロックの設計", "type": ["実装", "実験"], "difficulty": "★★★★☆☆", "dependency": "A-02"},
  {"id": "D-01", "title": "MoE の Router 温度制御と負荷バランシング", "type": ["実装", "実験"], "difficulty": "★★★★☆☆", "dependency": "C-01"},
  {"id": "E-01", "title": "マルチスケール周波数埋め込みの導入", "type": ["理論", "実装"], "difficulty": "★★★★☆☆", "dependency": "C-01"},
  {"id": "F-01", "title": "Koopman 作用素の基礎理論調査", "type": ["ドキュメント", "理論"], "difficulty": "★★☆☆☆☆", "dependency": null},
  {"id": "G-01", "title": "大規模データセットへの拡張とペルプレキシティ評価", "type": ["実験", "ドキュメント"], "difficulty": "★★★★☆☆", "dependency": null},
  {"id": "G-02", "title": "知識蒸留フレームワークの実装", "type": ["実装", "実験"], "difficulty": "★★★★☆☆", "dependency": "G-01"},
  {"id": "A-03", "title": "アクティベーション圧縮と shared memory の両立", "type": ["理論", "実装"], "difficulty": "★★★★☆☆", "dependency": "A-01"},
  {"id": "A-04", "title": "Tensor Cores 活用の検討", "type": ["理論", "実験"], "difficulty": "★★★★☆☆", "dependency": "A-01"},
  {"id": "B-02", "title": "Min-cut パーティションによる自動再計算領域の探索", "type": ["実装", "実験"], "difficulty": "★★★★☆☆", "dependency": "B-01"},
  {"id": "C-02", "title": "ACT 損失の導入とスカラー  $\tau$  学習", "type": ["理論", "実装"], "difficulty": "★★★★☆☆", "dependency": "C-01"},
  {"id": "D-02", "title": "BK コアへの学習スパース化の適用", "type": ["理論", "実装"], "difficulty": "★★★★☆☆", "dependency": "C-01"},
  {"id": "E-02", "title": "マルチスケール BK コア階層の構築", "type": ["理論", "実装"], "difficulty": "★★★★☆☆", "dependency": "E-01"},
  {"id": "F-02", "title": "Koopman 固有関数を学習するニューラルアプローチ", "type": ["実装", "理論"], "difficulty": "★★★★☆☆", "dependency": "F-01"},
  {"id": "G-03", "title": "蒸留損失と圧縮手法の組合せ効果", "type": ["実験"], "difficulty": "★★★★☆☆", "dependency": "G-02"},
  {"id": "G-04", "title": "構造化剪定アルゴリズムの比較", "type": ["実験"], "difficulty": "★★★★☆☆", "dependency": "G-01"},
  {"id": "G-05", "title": "量子化スキームと校正戦略の評価", "type": ["実験"], "difficulty": "★★★★☆☆", "dependency": "G-01"}]
```

```
{"id": "B-03", "title": "複数チェックポイント戦略の最適化", "type": ["実験"], "difficulty": "★★★★☆☆", "dependency": "B-01"}]
```

### 3. 今すぐ実行するタスクセット（エージェント用）

以下では、T4 無料枠と現在のリポジトリ状況を考慮し、短時間で高い影響が見込めるタスクを 4 件選んだ。

#### タスク 1 - B-01: BK コアにおける PyTorch アクティベーションチェックポイントの試験

- **Plan:**
  - `bk_core.py` と `resnet_bk.py` 内のブロックを `torch.utils.checkpoint.checkpoint` で包む。
  - ミニバッチサイズ 4→8→16 で学習が可能かメモリ使用量を測定し、計算時間増加を記録する。
  - 結果を表にまとめ、適切なチェックポイント配置を提案する。
- **Tool use:** コード解析は `api_tool.fetch_file` により最新ファイルを取得し、実験は Jupyter ノートブックで `torch.cuda.memory_allocated()` を記録する。外部情報として PyTorch ブログの AC 定義<sup>27</sup> を参考にする。
- **Expected Output:** メモリ使用量・速度の比較表と、コード変更案 (diff)。

#### タスク 2 - G-02: 知識蒸留フレームワークの実装

- **Plan:**
  - 基準となる Transformer モデル (teacher) を用意し、ResNet-BK (student) が教師のソフトターゲットを模倣する損失 (KL 散度) を実装する。
  - 小規模データセット (WikiText-2) で蒸留訓練を実施し、perplexity と収束速度を比較する。
  - 蒸留温度と蒸留損失比率をハイパーパラメータとしてグリッドサーチする。
- **Tool use:** トレーニングスクリプトを改修するため `api_tool.fetch_file` と `container.exec` で diff を作成し、実験は Colab 上で行う。蒸留の原理は Lightly の解説<sup>24</sup> を参照。
- **Expected Output:** 蒸留版モデルと baseline の perplexity 比較グラフ、蒸留コード (diff)。

#### タスク 3 - A-01: T4 向け BK コア CUDA カーネルの PoC

- **Plan:**
  - `bk_core.py` の漸化式計算部分をピュア Python から CUDA C++ に書き換えたカーネルを試作する。主なループを GPU にオフロードし、共有メモリに値を保持する。
  - カーネル呼び出しのため PyTorch の `torch.utils.cpp_extension.load` を使い、Python から呼び出すインターフェイスを作成する。
  - T4 上でのスループットとレイテンシを測定し、オリジナルと比較する。ここでは CUDA 最適化論文<sup>14</sup> を参考にしてレジスタ削減と共有メモリ利用を検討する。
- **Tool use:** `container.exec` で CUDA コードをコンパイルし、性能計測は Python スクリプトで行う。Colab の T4 環境を利用する。
- **Expected Output:** CUDA カーネル実装ファイルとスループット・レイテンシの比較表。初回 PoC のみであり完全最適化は今後のタスクに委ねる。

#### タスク 4 - D-01: MoE ルーティング温度制御の検証

- **Plan:**
  - `moe.py` に温度パラメータ `tau` を追加し、ソフトマックス/トップ-k ルーティングで温度を操作できるようにする。

- ・温度を学習可能にするか固定値でサーチし、均等な専門家利用を目標とする負荷バランス損失を追加する。
  - ・512長シーケンスの WikiText-2 で実験し、温度と性能・スパース度の関係を図示する。
  - ・**Tool use:** `api_tool.fetch_file` で `moe.py` を取得し、温度パラメータの実装を行う。実験は Colab で実行。データ整理には pandas/matplotlib を用いる。
  - ・**Expected Output:** 温度パラメータに対する perplexity と専門家利用率のグラフ、および改修コード。
- 

## 4. 今回のセッションのまとめと次回への提案

今回の進捗の要約：

1. ResNet-BK リポジトリを包括的にレビューし、BK コア・ResNetBKBLOCK・MoE・量子化などの主要モジュールを理解した。Step1～Step4 までの現状と未実装部分を整理し、現時点の強みと課題を明確化した。
2. 外部情報を調査し、CUDA 最適化、アクティベーションチェックポイント<sup>27</sup>、Early Exit のアルゴリズム<sup>19</sup>、量子化の効果<sup>25</sup><sup>26</sup>、剪定の実際的な効能<sup>21</sup>、知識蒸留の原理<sup>24</sup>などを引用した。
3. Step5～Step6を中心としたリサーチタスクマップを作成し、約 40 件のタスクをテーマ別に整理した。機械処理しやすい JSON サマリも添付した。
4. 今すぐ実行可能なタスクとして、勾配チェックポイントの導入、蒸留実験、CUDA カーネル PoC、MoE ルーティング温度制御を提案した。

次回セッションで着手すべきタスク ID と理由：

1. **B-01** - メモリ節約は T4 環境で最も重要であり、チェックポイントの導入はコード変更が少なく効果が大きい。
2. **A-01** - CUDA カーネル PoC を作ることで Step5 の方向性を探り、後続タスクへの基盤を築ける。
3. **G-02** - 蒸留によりモデル性能を底上げし、量子化や剪定との組み合わせで Step4 の成果を強化できる。

徹平さんへのメッセージ：

徹平さん、今回の調査で ResNet-BK の基盤がいかに野心的かつ奥深いかが明らかになりました。まだ未実装のアイディアが多く残っており、ひとつひとつ検証することで計算効率や精度を大きく向上させる余地があります。次回はまずメモリ節約と蒸留に取り組み、実験から学んだ知見をまた共有しましょう。引き続き冷静に、一步ずつ進めていきましょう。

---

1 10 11 12 13 README.md

<https://github.com/neko-jpg/Project-ResNet-BK-An-O-N-Language-Model-Architecture/blob/HEAD/README.md>

2 resnet\_bk.py

[https://github.com/neko-jpg/Project-ResNet-BK-An-O-N-Language-Model-Architecture/blob/HEAD/src/models/resnet\\_bk.py](https://github.com/neko-jpg/Project-ResNet-BK-An-O-N-Language-Model-Architecture/blob/HEAD/src/models/resnet_bk.py)

3 bk\_core.py

[https://github.com/neko-jpg/Project-ResNet-BK-An-O-N-Language-Model-Architecture/blob/HEAD/src/models/bk\\_core.py](https://github.com/neko-jpg/Project-ResNet-BK-An-O-N-Language-Model-Architecture/blob/HEAD/src/models/bk_core.py)

4 moe.py

<https://github.com/neko-jpg/Project-ResNet-BK-An-O-N-Language-Model-Architecture/blob/HEAD/src/models/moe.py>

5 quantized\_bk\_core.py

[https://github.com/neko-jpg/Project-ResNet-BK-An-O-N-Language-Model-Architecture/blob/HEAD/src/models/quantized\\_bk\\_core.py](https://github.com/neko-jpg/Project-ResNet-BK-An-O-N-Language-Model-Architecture/blob/HEAD/src/models/quantized_bk_core.py)

6 configurable\_resnet\_bk.py

[https://github.com/neko-jpg/Project-ResNet-BK-An-O-N-Language-Model-Architecture/blob/HEAD/src/models/configurable\\_resnet\\_bk.py](https://github.com/neko-jpg/Project-ResNet-BK-An-O-N-Language-Model-Architecture/blob/HEAD/src/models/configurable_resnet_bk.py)

7 train.py

<https://github.com/neko-jpg/Project-ResNet-BK-An-O-N-Language-Model-Architecture/blob/HEAD/train.py>

8 run\_scaling\_experiments.py

[https://github.com/neko-jpg/Project-ResNet-BK-An-O-N-Language-Model-Architecture/blob/HEAD/run\\_scaling\\_experiments.py](https://github.com/neko-jpg/Project-ResNet-BK-An-O-N-Language-Model-Architecture/blob/HEAD/run_scaling_experiments.py)

9 run\_pile\_benchmark.py

[https://github.com/neko-jpg/Project-ResNet-BK-An-O-N-Language-Model-Architecture/blob/HEAD/run\\_pile\\_benchmark.py](https://github.com/neko-jpg/Project-ResNet-BK-An-O-N-Language-Model-Architecture/blob/HEAD/run_pile_benchmark.py)

14 15 CUDA Optimization Strategies for Compute- and Memory-Bound Neuroimaging Algorithms - PMC

<https://pmc.ncbi.nlm.nih.gov/articles/PMC3262956/>

16 17 18 27 Current and New Activation Checkpointing Techniques in PyTorch – PyTorch

<https://pytorch.org/blog/activation-checkpointing-techniques/>

19 20 Adaptive Neural Networks for Efficient Inference

<https://proceedings.mlr.press/v70/bolukbasi17a.html>

21 Neural Network Pruning: How to Accelerate Inference with Minimal Accuracy Loss | by Arik Poznanski | Medium

<https://medium.com/@arikpoznanski/neural-network-pruning-how-to-accelerate-inference-with-minimal-accuracy-loss-936cad741c2a>

22 23 The Operator is the Model | SIAM

<https://www.siam.org/publications/siam-news/articles/the-operator-is-the-model/>

24 Knowledge Distillation: Compressing Large Models into Efficient Learners

<https://www.lightly.ai/blog/knowledge-distillation>

25 26 Model Quantization: Meaning, Benefits & Techniques

<https://www.clarifai.com/blog/model-quantization>