

CASPR: Combining Axis Preconditioners through Kronecker Approximation for Deep Learning

Sai Surya Duvvuri

Joint work with Devvrit, Rohan Anil, Cho-Jui Hsieh, Inderjit Dhillon



Optimization Algorithms in Practice

- Objective:

$$\min_w f(w), \quad w \in \mathbb{R}^d$$

- Solution:

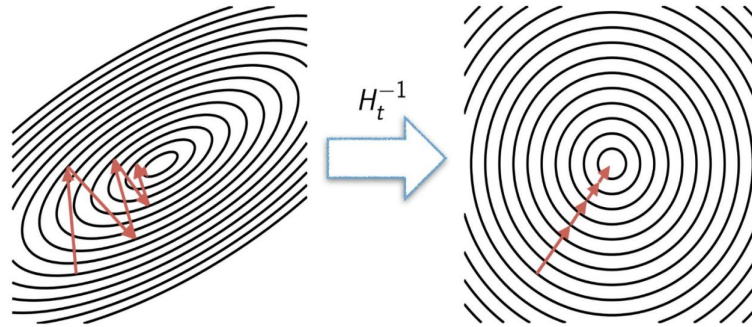
for $t = 2, \dots, T$

$$w_t = w_{t-1} - \eta_t X_t \nabla f(w_{t-1})$$

- X_t is a preconditioner matrix of size $d \times d$.

Preconditioning

often leads to faster convergence / better "condition number"



Geometrically they scale and rotate gradients

- Preconditioning typically involves inverting curvature information.

Optimizers for Large Deep Learning Models

- Finding preconditioner can incur high memory and compute.
- diagonal preconditioners:
 - Adam and Adagrad use coordinate wise second-moments $(g)_i^2$
 - But don't utilize cross moments $(g)_i(g)_j$
- Full-matrix Adagrad uses cross-moments → potential for faster convergence!

Aim

- Full matrix-Adagrad update:

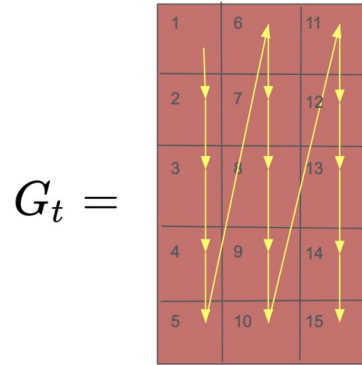
$$H_t = H_{t-1} + g_t g_t^T, \quad \rightarrow \text{memory intensive} - \mathcal{O}(d^2)$$
$$w_{t+1} = w_t - \eta H_t^{-1/2} g_t \quad \rightarrow \text{compute intensive} - \mathcal{O}(d^3)$$

- Develop an approximation \hat{H}_t to second-moment matrix H_t :
 - Accurate approximation.
 - Low memory to store.
 - Fast inversion.

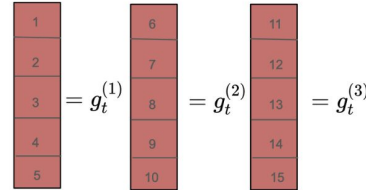
Shampoo (Gupta et al., 2018)

- Scalable implementation (Anil et al., 2020)
- Applications in recommendation models in Google (Anil et al., 2022)
- Practical kronecker product approximation.
- We utilize Kronecker sum to develop a better approximation.

Approximating Second-Moment Matrix of 2-D Parameter

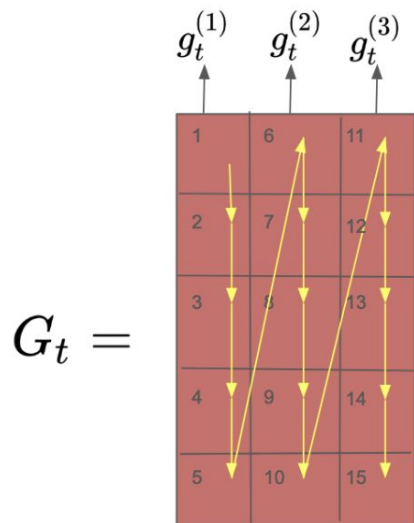


Column Major Flattening



Each column as separate gradient vector

Block-Diagonal Approximation with Identical Blocks



Column Major Flattening

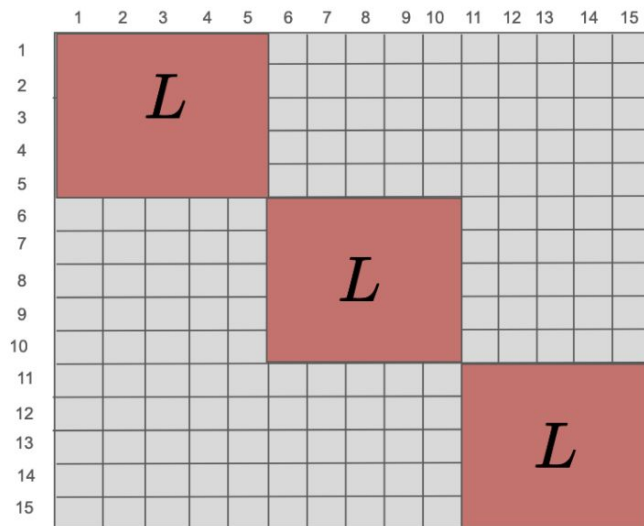


Figure: $I_n \otimes L$

Block-Diagonal Approximation with Identical Blocks

- Set all blocks to be L in the subproblem

$$L^* = \operatorname{argmin}_{L \succeq 0} \left\| I_n \otimes L - \sum_{t=1}^T g_t g_t^\top \right\|_F = \frac{1}{n} \sum_{t=1}^T G_t G_t^\top \quad (\text{Explicit Solution!})$$

Row Preconditioner

- We can similarly form row-preconditioner.

$$R^* = \operatorname{argmin}_{L \succeq 0} \left\| R \otimes I_m - \sum_{t=1}^T g_t g_t^\top \right\|_F = \frac{1}{m} \sum_{t=1}^T G_t^\top G_t$$

Axes Preconditioners

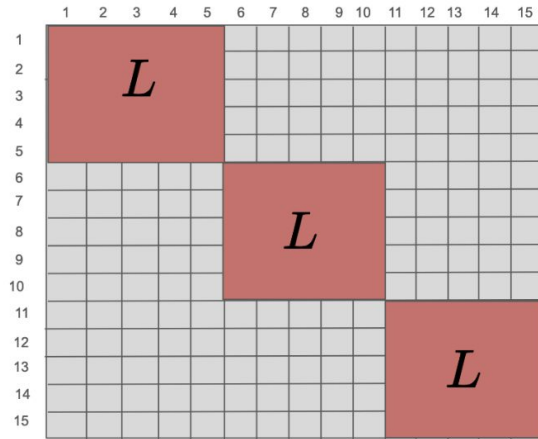


Figure: $I_n \otimes L$

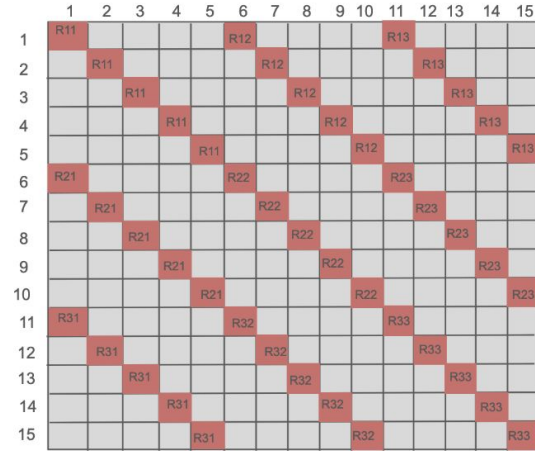


Figure: $R \otimes I_m$

- Individually both approximations miss out on a lot of cross-moments.

Axes Preconditioners

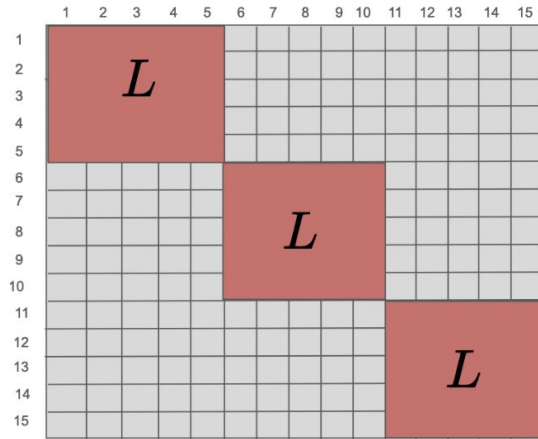


Figure: $I_n \otimes L$

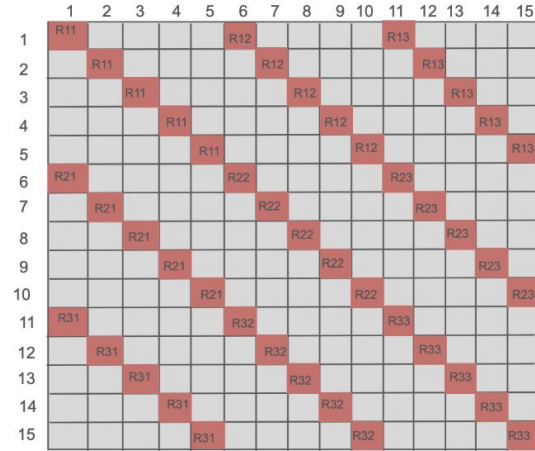


Figure: $R \otimes I_m$

- Individually both approximations miss out on a lot of cross-moments.
- Should combine both to approximate the remaining cross-moments?

CASPR Update

- CASPR update for $p = 2$ is:

$$X_t := \left((R_t^{-1/4} \otimes I_m + I_n \otimes L_t^{-1/4}) / 2 \right)^2; \quad W_t := W_{t-1} - \eta X_t g_t,$$

- Expanding the update gives:

$$W_{t+1} := W_t - \eta \left(L_t^{-1/2} G_t + 2L_t^{-1/4} G_t R_t^{-1/4} + G_t R_t^{-1/2} \right)$$

Comparison with Shampoo Update

CASPR update

Shampoo update

Update preconditioners:

$$L_t := L_{t-1} + G_t G_t^\top, \quad R_t := R_{t-1} + G_t^\top G_t$$

Compute $L^{-1/4}$, $R^{-1/4}$

Precondition gradient and update parameters:

$$U_t := L_t^{-1/4} G_t + G_t R_t^{-1/4}$$

$$U_t := L_t^{-1/4} U_t + U_t R_t^{-1/4}$$

$$W_t := W_{t-1} - \eta U_t$$

$$U_t := L_t^{-1/4} G_t R_t^{-1/4}$$

$$W_t := W_{t-1} - \eta U_t$$

Regret Bound Analysis

- We conduct analysis in online convex optimization framework.

Theorem (Regret upper bound of CASPR Algorithm)

Given that the loss functions $f_t : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$, $\forall t \in [T]$ are convex and G -Lipschitz in ℓ_2 -norm i.e., $\|\nabla f_t(W)\|_2 \leq G$, $W \in \mathbb{R}^{m \times n}$, Algorithm 1 incurs the following regret

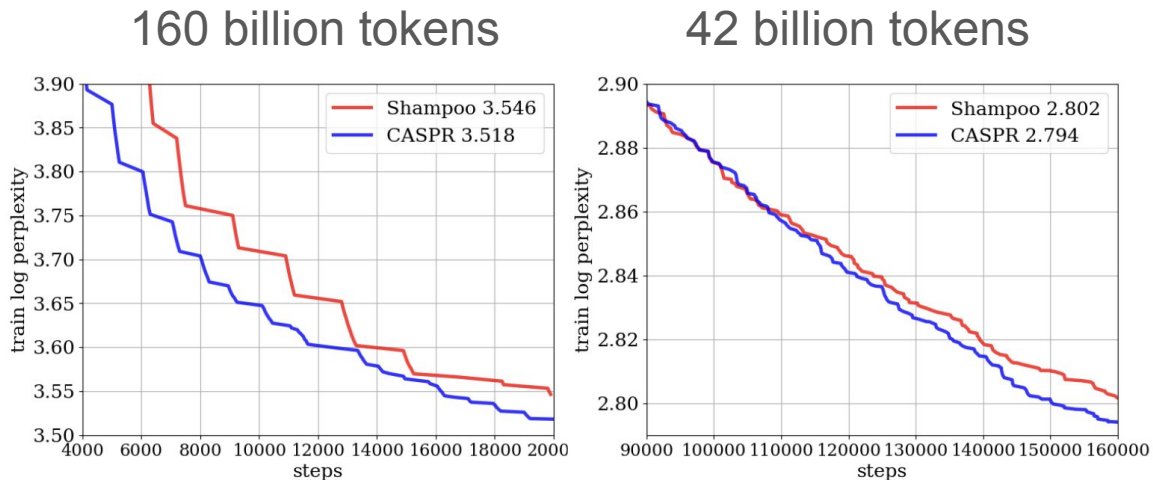
$$\begin{aligned} \sum_{t=1}^T f_t(W_t) - f_t(W^*) &\leq \sqrt{2r}D \operatorname{tr} \left(\left((L_T^{-1/4} \otimes I_n + I_m \otimes R_T^{-1/4})/2 \right)^{-2} \right) \\ &\leq \sqrt{2r}D \operatorname{tr} \left(L_T^{1/4} \otimes R_T^{1/4} \right) = \mathcal{O}(\sqrt{T}) \end{aligned}$$

when $\eta = D/\sqrt{2r}$, where $r = \max_t \operatorname{rank}(G_t)$, $D = \|W_t - W^*\|_F$

- CASPR has tighter regret upper bound than Shampoo.

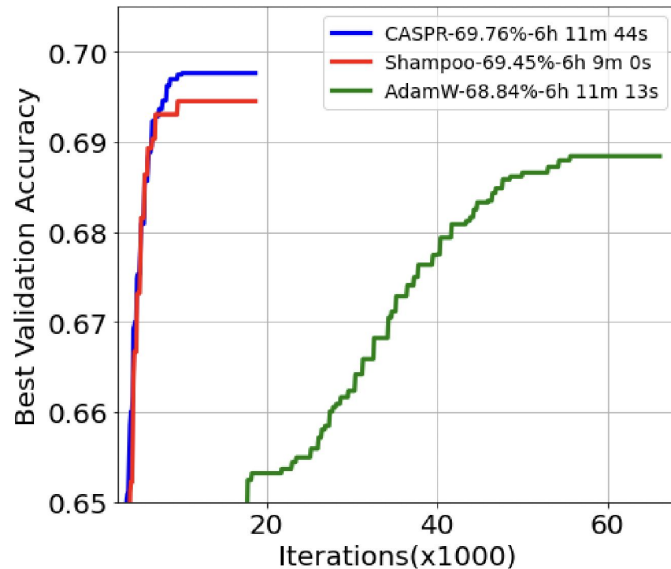
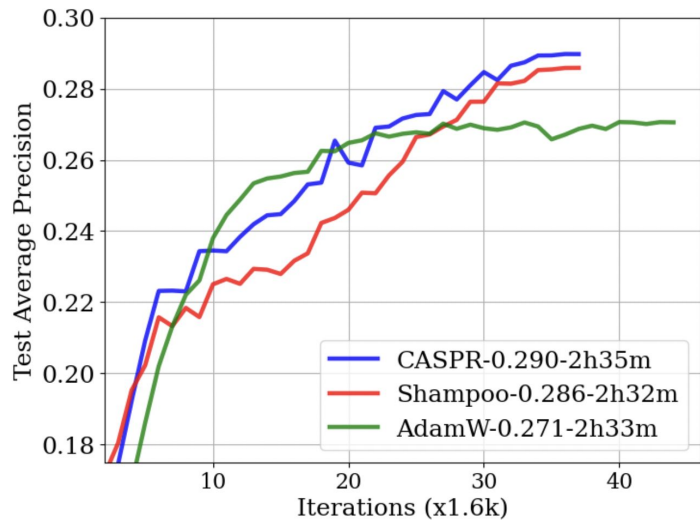
Autoregressive Large Language Modeling

- GLU based decoder-only transformer models trained on C4 dataset.



a) 8192 batch size and 14M parameters b) 256 batch size and 234M parameters.

GNN and Transformer on Parts of Speech



- Time taken for Shampoo and CASPR are about the same.
- CASPR has a better Validation Accuracy than Shampoo.
- CASPR is better than AdamW when run for fixed amount of time.

Conclusion and Future Directions

- Novel Kronecker-sum inspired combination approach to approximate the second-moment matrix using axes preconditioners.
- Stronger convergence guarantees than Shampoo, which is a special case of our framework of combining axes preconditioners.
- More accurate axes preconditioners solving the problem

$$\hat{H}_t := \arg \min_{\hat{H} \in \mathcal{S}} \left\| \hat{H} - H_t \right\|_F$$

- Adapt CASPR to approximate Hessian instead of full-matrix Adagrad.