

## 490 Appendix Outline

491 This Appendix is organized as follows:

- 492 • In [Appendix A](#) we describe various dispatch rules including the base rules, the composition
- 493 rules and rules derived from other rules.
- 494 • In [Appendix B](#) we provide an extended discussion of several noteworthy features of CoLA,
- 495 such as doubly stochastic estimators and memory-efficient autograd implementation.
- 496 • In [Appendix C](#) we include pseudo-code on various of the iterative methods incorporated in
- 497 CoLA and discuss modifications to improve lower precision performance.
- 498 • In [Appendix D](#) we expand on the details of the experiments in the main text.
- 499 • In [Appendix E](#) we show some code examples of how the dispatch rules are implemented in
- 500 CoLA.

## 501 A Dispatch Rules

502 We now present the linear algebra identities that we use to exploit structure in CoLA.

### 503 A.1 Core Functions

#### 504 A.1.1 Inverses

505 We incorporate several identities for the compositional operators: product, Kronecker product, block  
 506 diagonal and sum. For product we have  $(\mathbf{AB})^{-1} = (\mathbf{B}^{-1}\mathbf{A}^{-1})$  and for Kronecker product we have  
 507  $(\mathbf{A} \otimes \mathbf{B})^{-1} = \mathbf{A}^{-1} \otimes \mathbf{B}^{-1}$ . In terms of block compositions we have the following identities:

$$\begin{aligned}
 & \begin{bmatrix} \mathbf{A} & \mathbf{0} \\ \mathbf{0} & \mathbf{D} \end{bmatrix}^{-1} = \begin{bmatrix} \mathbf{A}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{D}^{-1} \end{bmatrix} \quad \text{and} \quad \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{0} & \mathbf{D} \end{bmatrix}^{-1} = \begin{bmatrix} \mathbf{A}^{-1} & -\mathbf{A}^{-1}\mathbf{B}\mathbf{D}^{-1} \\ \mathbf{0} & \mathbf{D}^{-1} \end{bmatrix} \\
 & \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}^{-1} = \begin{bmatrix} \mathbf{I} & -\mathbf{A}^{-1}\mathbf{B} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{A} & \mathbf{0} \\ \mathbf{0} & \mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ -\mathbf{C}\mathbf{A}^{-1} & \mathbf{I} \end{bmatrix}
 \end{aligned}$$

509 Finally, for sum we have the Woodbury identity and its variants. Namely, for Woodbury we have

$$(\mathbf{A} + \mathbf{UBV})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{U}(\mathbf{B}^{-1} + \mathbf{VA}^{-1}\mathbf{U})^{-1}\mathbf{VA}^{-1},$$

510 the Kailath variant where

$$(\mathbf{A} + \mathbf{BC})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{B}(\mathbf{I} + \mathbf{C}\mathbf{A}^{-1}\mathbf{B})\mathbf{C}\mathbf{A}^{-1}$$

511 and the rank one update via the Sherman-Morrison formula

$$(\mathbf{A} + \mathbf{bc}^\top)^{-1} = \mathbf{A}^{-1} - \frac{1}{1 + \mathbf{c}^\top\mathbf{A}\mathbf{b}}\mathbf{A}^{-1}\mathbf{bc}^\top\mathbf{A}^{-1}.$$

512 Besides the compositional operators, we have some rules for some special operators. For example,  
 513 for  $\mathbf{A} = \text{Diag}(\mathbf{a})$  we have  $\mathbf{A}^{-1} = \text{Diag}(\mathbf{a}^{-1})$ . Also, if  $\mathbf{Q}$  is unitary then  $\mathbf{Q}^{-1} = \mathbf{Q}^*$  or if  $\mathbf{Q}$  is  
 514 orthonormal then  $\mathbf{Q}^{-1} = \mathbf{Q}^\top$ . In [Appendix E](#) we show how these dispatch rules are implemented in  
 515 Python.

#### 516 A.1.2 Eigendecomposition

517 We now assume that the matrices in this section are diagonalizable. That is,  $\text{Eigs}(\mathbf{A}) = \Lambda_{\mathbf{A}}, \mathbf{V}_{\mathbf{A}}$ ,  
 518 where  $\mathbf{A} = \mathbf{V}_{\mathbf{A}}\Lambda_{\mathbf{A}}\mathbf{V}_{\mathbf{A}}^{-1}$ . In terms of the compositional operators, there is not a general rule for  
 519 product or sum. However, for the Kronecker product we have  $\text{Eigs}(\mathbf{A} \otimes \mathbf{B}) = \Lambda_{\mathbf{A}} \otimes \Lambda_{\mathbf{B}}, \mathbf{V}_{\mathbf{A}} \otimes \mathbf{V}_{\mathbf{B}}$   
 520 and for the Kronecker sum we have  $\text{Eigs}(\mathbf{A} \oplus \mathbf{B}) = \Lambda_{\mathbf{A}} \oplus \Lambda_{\mathbf{B}}, \mathbf{V}_{\mathbf{A}} \otimes \mathbf{V}_{\mathbf{B}}$ . Finally, for block  
 521 diagonal we have

$$\text{Eigs}\left(\begin{bmatrix} \mathbf{A} & \mathbf{0} \\ \mathbf{0} & \mathbf{D} \end{bmatrix}\right) = \begin{bmatrix} \Lambda_{\mathbf{A}} & \mathbf{0} \\ \mathbf{0} & \Lambda_{\mathbf{D}} \end{bmatrix}, \begin{bmatrix} \mathbf{V}_{\mathbf{A}} & \mathbf{0} \\ \mathbf{0} & \mathbf{V}_{\mathbf{D}} \end{bmatrix}.$$

### 522 A.1.3 Diagonal

523 As a base case, if we need to compute  $\text{Diag}(\mathbf{A})$  for a general matrix  $\mathbf{A}$  we may compute each  
 524 diagonal element by  $\mathbf{e}_i^\top \mathbf{A} \mathbf{e}_i$ . Additionally, if  $\mathbf{A}$  is large enough we switch to randomized estimation  
 525  $\text{Diag}(\mathbf{A}) \approx (\mathbf{Z} \odot \mathbf{A} \mathbf{Z}) \mathbf{1} / N$  with  $\mathbf{Z} \sim \mathcal{N}(0, 1)^{d \times N}$  where  $N$  is the number of samples used to ap-  
 526 proximate the diagonal. In terms of compositional operators, we have that for sum  $\text{Diag}(\mathbf{A} + \mathbf{B}) =$   
 527  $\text{Diag}(\mathbf{A}) + \text{Diag}(\mathbf{B})$ . For Kronecker product we have  $\text{Diag}(\mathbf{A} \otimes \mathbf{B}) = \text{vec}(\text{Diag}(\mathbf{A}) \text{Diag}(\mathbf{B})^\top)$   
 528 and for Kronecker sum  $\text{Diag}(\mathbf{A} \oplus \mathbf{B}) = \text{vec}(\text{Diag}(\mathbf{A}) \mathbf{1}^\top + \mathbf{1} \text{Diag}(\mathbf{B})^\top)$ . Finally, for block  
 529 composition we have

$$\text{Diag}\left(\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}\right) = [\text{Diag}(\mathbf{A}), \text{Diag}(\mathbf{D})].$$

### 530 A.1.4 Transpose / Adjoint

531 As explained in [Section 3.1](#), as a base case we have an automatic procedure to compute the transpose or  
 532 adjoint of any operator  $\mathbf{A}$  via autodiff. However, we also incorporate the following rules. For sum we  
 533 have  $(\mathbf{A} + \mathbf{B})^* = \mathbf{A}^* + \mathbf{B}^*$  and  $(\mathbf{A} + \mathbf{B})^\top = \mathbf{A}^\top + \mathbf{B}^\top$ . For product we have  $(\mathbf{A}\mathbf{B})^* = \mathbf{B}^* \mathbf{A}^*$  and  
 534  $(\mathbf{A}\mathbf{B})^\top = \mathbf{B}^\top \mathbf{A}^\top$ . For Kronecker product we have  $(\mathbf{A} \otimes \mathbf{B})^* = \mathbf{A}^* \otimes \mathbf{B}^*$  and  $(\mathbf{A} \otimes \mathbf{B})^\top = \mathbf{A}^\top \otimes \mathbf{B}^\top$ .  
 535 For the Kronecker sum we have  $(\mathbf{A} \oplus \mathbf{B})^* = \mathbf{A}^* \oplus \mathbf{B}^*$  and  $(\mathbf{A} \oplus \mathbf{B})^\top = \mathbf{A}^\top \oplus \mathbf{B}^\top$ . In terms of  
 536 block composition we have

$$\left(\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}\right)^* = \begin{bmatrix} \mathbf{A}^* & \mathbf{C}^* \\ \mathbf{B}^* & \mathbf{D}^* \end{bmatrix} \quad \text{and} \quad \left(\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}\right)^\top = \begin{bmatrix} \mathbf{A}^\top & \mathbf{C}^\top \\ \mathbf{B}^\top & \mathbf{D}^\top \end{bmatrix}.$$

537 Finally for the annotated operators we have the following rules.  $\mathbf{A}^* = \mathbf{A}$  if  $\mathbf{A}$  is self-adjoint and  
 538  $\mathbf{A}^\top = \mathbf{A}$  if  $\mathbf{A}$  is symmetric.

### 539 A.1.5 Pseudo-inverse

540 As a base case, if we need to compute  $\mathbf{A}^+$ , we may use  $\text{SVD}(\mathbf{A}) = \mathbf{U}, \Sigma, \mathbf{V}$  and therefore set  
 541  $\mathbf{A}^+ = \mathbf{U} \Sigma^+ \mathbf{V}^*$ , where  $\Sigma^+$  inverts the nonzero diagonal scalars. If the size of  $\mathbf{A}$  is too large,  
 542 then we may use randomized SVD. Yet, it is uncommon to simply want  $\mathbf{A}^+$ , usually we want to  
 543 solve a least-squares problem and therefore we can use solvers that are not as expensive to run as  
 544 SVD. For the compositional operators we have the following identities. For product  $(\mathbf{A}\mathbf{B})^+ =$   
 545  $(\mathbf{A}^+ \mathbf{A} \mathbf{B})^+ (\mathbf{A} \mathbf{B} \mathbf{B}^+)^+$  and for Kronecker product we have  $(\mathbf{A} \otimes \mathbf{B})^+ = \mathbf{A}^+ \otimes \mathbf{B}^+$ . For block  
 546 diagonal we have

$$\left(\begin{bmatrix} \mathbf{A} & \mathbf{0} \\ \mathbf{0} & \mathbf{D} \end{bmatrix}\right)^+ = \begin{bmatrix} \mathbf{A}^+ & \mathbf{0} \\ \mathbf{0} & \mathbf{D}^+ \end{bmatrix}.$$

547 Finally, we have some identities that are mathematically trivial but that are necessary when recursively  
 548 exploiting structure as that would save computation. For example, if  $\mathbf{Q}$  is unitary we know that  
 549  $\mathbf{Q}^+ = \mathbf{Q}$  and similarly when  $\mathbf{Q}$  is orthonormal. If  $\mathbf{A}$  is self-adjoint, then  $\mathbf{A}^+ = \mathbf{A}^{-1}$  and also if it is  
 550 symmetric and PSD.

## 551 A.2 Derived Functions

552 Interestingly, the previous core functions allow us to derive multiple rules from the previous ones.  
 553 To illustrate, we have that  $\text{Tr}(\mathbf{A}) = \sum_i \text{Diag}(\mathbf{A})_i$ . Additionally, if  $\mathbf{A}$  is PSD we have that  
 554  $f(\mathbf{A}) = \mathbf{V}_\mathbf{A} f(\Lambda_\mathbf{A}) \mathbf{V}_\mathbf{A}^{-1}$  and if  $\mathbf{A}$  is both symmetric and PSD then  $f(\mathbf{A}) = \mathbf{V}_\mathbf{A} f(\Lambda_\mathbf{A}) \mathbf{V}_\mathbf{A}^\top$ .  
 555 where in both cases we used  $\text{Eigs}(\mathbf{A}) = \Lambda_\mathbf{A}, \mathbf{V}_\mathbf{A}$ . Some example functions for PSD matrices are  
 556  $\text{Sqrt}(\mathbf{A}) = \mathbf{V}_\mathbf{A} \Lambda_\mathbf{A}^{1/2} \mathbf{V}_\mathbf{A}^{-1}$  or  $\text{Log}(\mathbf{A}) = \mathbf{V}_\mathbf{A} \log \Lambda_\mathbf{A} \mathbf{V}_\mathbf{A}^{-1}$ . Which also this rules allow us to define  
 557  $\text{LogDet}(\mathbf{A}) = \text{Tr}(\text{Log}(\mathbf{A}))$ .

## 558 A.3 Other matrix identities

559 We emphasize that there are a myriad more matrix identities that we do not intentionally include  
 560 such as  $\text{Tr}(\mathbf{A} + \mathbf{B}) = \text{Tr}(\mathbf{A}) + \text{Tr}(\mathbf{B})$  or  $\text{Tr}(\mathbf{A}\mathbf{B}) = \text{Tr}(\mathbf{B}\mathbf{A})$  when  $\mathbf{A}$  and  $\mathbf{B}$  are squared. These

561 additional cases are not part of our dispatch rules as either they are automatically computed from  
562 other rules (as in the first example) or they do not yield any computational savings (as in the second  
563 example).

## 564 B Features in CoLA

### 565 B.1 Doubly stochastic diagonal and trace estimation

566 **Singly Stochastic Trace Estimator** Consider the traditional stochastic trace estimator:

$$\overline{\text{Tr}}[\text{Base}](\mathbf{A}) = \frac{1}{n} \sum_{j=1}^n \mathbf{z}_j^\top \mathbf{A} \mathbf{z}_j \quad (1)$$

567 with each  $\mathbf{z}_j \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_D)$  where  $\mathbf{A}$  is a  $D \times D$  matrix. When  $\mathbf{A}$  is itself a sum  $\mathbf{A} = \frac{1}{m} \sum_{i=1}^m \mathbf{A}_i$ ,  
568 we can expand the trace as  $\overline{\text{Tr}}[\text{Base}](\mathbf{A}) = \frac{1}{mn} \sum_{j=1}^n \sum_{i=1}^m \mathbf{z}_j^\top \mathbf{A}_i \mathbf{z}_j$ , with probe variables shared  
569 across elements of the sum.

570 Consider the quadratic form  $Q := \mathbf{z}^\top \mathbf{A} \mathbf{z}$ , which for Gaussian random variables has a cumulant  
571 generating function of  $K_Q(t) = \log \mathbb{E}[e^{tQ}] = -\frac{1}{2} \log \det(\mathbf{I} - 2t\mathbf{A})$ . From the generating function  
572 we can derive the mean and variance of this estimator:  $\mathbb{E}[Q] = K'_Q(0) = \text{Tr}(\mathbf{A})$  and  $\text{Var}[Q] =$   
573  $K''_Q(0) = 2\text{Tr}(\mathbf{A}^2)$ . Since  $\overline{\text{Tr}}[\text{Base}](\mathbf{A})$  is a sum of independent random draws of  $Q$ , we see:

$$\mathbb{E}[\overline{\text{Tr}}[\text{Base}](\mathbf{A})] = \text{Tr}(\mathbf{A}) \quad \text{and} \quad \text{Var}[\overline{\text{Tr}}[\text{Base}](\mathbf{A})] = \frac{2}{n} \text{Tr}(\mathbf{A}^2). \quad (2)$$

574 **Doubly Stochastic Trace Estimator** For the doubly stochastic estimator, we choose probe vari-  
575 ables which are sampled independently for each element of the sum:

$$\overline{\text{Tr}}[\text{Sum}](\mathbf{A}) = \frac{1}{nm} \sum_{j=1}^n \sum_{i=1}^m \mathbf{z}_{ij}^\top \mathbf{A}_i \mathbf{z}_{ij}. \quad (3)$$

576 Separating out the elements of the sum, we can write the estimator as  $\overline{\text{Tr}}[\text{Sum}](\mathbf{A}) = \frac{1}{n} \sum_{j=1}^n R_j$   
577 where  $R_j$  are independent random samples of the value  $R = \frac{1}{m} \sum_{i=1}^m \mathbf{z}_i^\top \mathbf{A}_i \mathbf{z}_i$ . The cumulant  
578 generating function is merely  $K_R(t) = \sum_{i=1}^m K_{Q_i}(t/m)$  where  $Q_i = \mathbf{z}^\top \mathbf{A}_i \mathbf{z}$ . Taking derivatives  
579 we find that,

$$\mathbb{E}[R] = K'_R(0) = \frac{1}{m} \sum_{i=1}^m \text{Tr}(\mathbf{A}_i) = \text{Tr}(\mathbf{A}), \quad (4)$$

580

$$\text{Var}[R] = K''_R(0) = \frac{1}{m^2} \sum_{i=1}^m 2\text{Tr}(\mathbf{A}_i^2) = \frac{2}{m} \text{Tr}\left(\frac{1}{m} \sum_{i=1}^m \mathbf{A}_i^2\right) \quad (5)$$

581 Assuming bounded moments on  $\mathbf{A}_i$ , then both  $\mathbf{A} = \frac{1}{m} \sum_i \mathbf{A}_i$  and  $S(\mathbf{A}) = \frac{1}{m} \sum_i \mathbf{A}_i^2$  will converge  
582 to fixed values as  $m \rightarrow \infty$ . Given that  $\overline{\text{Tr}}[\text{Sum}](\mathbf{A}) = \frac{1}{n} \sum_{j=1}^n R_j$ , we can now write the mean and  
583 variance of the doubly stochastic estimator:

$$\mathbb{E}[\overline{\text{Tr}}[\text{Sum}](\mathbf{A})] = \text{Tr}(\mathbf{A}) \quad \text{and} \quad \text{Var}[\overline{\text{Tr}}[\text{Sum}](\mathbf{A})] = \frac{2}{mn} \text{Tr}(S(\mathbf{A})). \quad (6)$$

584 As the error of the estimator can be bounded by the square root of the variance, showing that while  
585 the error for  $\overline{\text{Tr}}[\text{Base}]$  is  $O(1/\sqrt{n})$  (even when applied to sum structures), whereas the error for  
586  $\overline{\text{Tr}}[\text{Sum}]$  is  $O(1/\sqrt{nm})$ , a significant asymptotic variance reduction.

587 The related stochastic diagonal estimator

$$\overline{\text{Diag}}[\text{Sum}](\mathbf{A}) = \frac{1}{nm} \sum_{j=1}^n \sum_{i=1}^m \mathbf{z}_{ij} \odot \mathbf{A}_i \mathbf{z}_{ij}. \quad (7)$$

588 achieves the same  $O(1/\sqrt{nm})$  convergence rate, though we omit this derivation for brevity as it is  
589 follows the same steps.

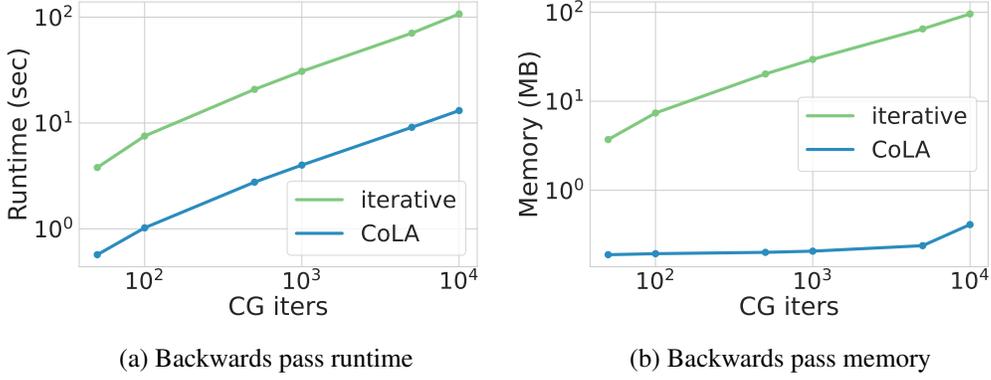


Figure 4: **Our autograd rules allow for fast and memory efficient backpropagation.** (a) Runtime required to compute  $\partial_{\theta} \mathbf{A}_{\theta}^{-1} \mathbf{b}$  for different solves, each requiring an increasing number of CG iterations. (a) Peak memory utilization required to compute  $\partial_{\theta} \mathbf{A}_{\theta}^{-1} \mathbf{b}$  for different solves, each requiring an increasing number of CG iterations.

## 590 B.2 Autograd rules for iterative algorithms

591 For machine learning applications, we want to seamlessly interweave linear algebra operations with  
 592 automatic differentiation. The most basic strategy is to simply let the autograd engine trace through  
 593 the operations and backpropagate accordingly. However, when using iterative methods like conjugate  
 594 gradients or Lanczos, this naive approach is extremely memory inefficient and, for problems with  
 595 many iterations, the cost can be prohibitive (as seen in Figure 4). However, the linear algebra  
 596 operations corresponding to inverse, eigendecomposition and trace estimation have simple closed  
 597 form derivatives which we can implement to avoid the prohibitive memory consumption and reduce  
 598 runtime.

599 Simply put, for an operation like  $f = \text{CGSolve}$ ,  $\text{CGSolve}(\mathbf{A}, \mathbf{b}) = \mathbf{A}^{-1} \mathbf{b}$  we must define a Vector  
 600 Jacobian Product:  $\text{VJP}(f, (\mathbf{A}, \mathbf{b}), \mathbf{v}) = (\mathbf{v}^{\top} \frac{\partial f}{\partial \mathbf{A}}, \mathbf{v}^{\top} \frac{\partial f}{\partial \mathbf{b}})$ . However, for matrix-free linear operators,  
 601 we cannot afford to store the dense matrix  $\mathbf{A}$ , and thus neither can we store the gradients with  
 602 respect to each of its elements! Instead we must (recursively) consider how the linear operator  
 603 was constructed in terms of its differentiable arguments. In other words, we must flatten the tree  
 604 structure of possibly nested differentiable arguments into a vector:  $\theta = \text{flatten}[\mathbf{A}]$ . For example  
 605 for  $\mathbf{A} = \text{Kron}(\text{Diag}(\theta_1), \text{Conv}(\theta_2))$ ,  $\text{flatten}[\mathbf{A}] = [\theta_1, \theta_2]$ . From this perspective, we consider  $\mathbf{A}$   
 606 as a container or tree of its arguments  $\theta$ , and define  $\mathbf{v}^{\top} \frac{\partial f}{\partial \mathbf{A}} := \text{unflatten}[\mathbf{v}^{\top} \frac{\partial f}{\partial \theta}]$  which coincides  
 607 with the usual definition for dense matrices. Applying to inverses, we can now write a simple VJP:

$$\mathbf{v}^{\top} \frac{\partial f}{\partial \mathbf{A}} = \text{unflatten}[\text{VJP}(\theta \mapsto \text{unflatten}(\theta) \mathbf{A}^{-1} \mathbf{b}, \theta, \mathbf{A}^{-1} \mathbf{v})] \quad (8)$$

608 for  $\mathbf{v}^{\top} \frac{\partial f}{\partial \theta} = \mathbf{v}^{\top} (\mathbf{A}^{-1})^{\top} (\partial_{\theta} \mathbf{A}_{\theta}) \mathbf{A}^{-1} \mathbf{b}$ , and we will adopt this notation below for brevity. Doing so  
 609 gives a memory cost which is constant in the number of solver iterations, and proportional to the  
 610 memory used in the forward pass. Below we list the autograd rules for some of the iterative routines  
 611 that we implement in CoLA with their VJP definitions.

- 612 1.  $\mathbf{y} = \text{Solve}(\mathbf{A}, \mathbf{b})$ :  $\mathbf{w}^{\top} \frac{\partial \mathbf{y}}{\partial \theta} = -(\mathbf{A}^{-1} \mathbf{w})^{\top} (\partial_{\theta} \mathbf{A}_{\theta}) (\mathbf{A}^{-1} \mathbf{b})$
- 613 2.  $\lambda, \mathbf{V} = \text{Eigs}(\mathbf{A})$ :  $\mathbf{w}^{\top} \frac{\partial \lambda}{\partial \theta} = \mathbf{w}^{\top} \text{Diag}(\mathbf{V}^{\top} (\partial_{\theta} \mathbf{A}_{\theta}) \mathbf{V})$
- 614 3.  $\lambda, \mathbf{V} = \text{Eigs}(\mathbf{A})$ :  $\mathbf{w}^{\top} \frac{\partial \mathbf{v}_i}{\partial \theta} = \mathbf{w}^{\top} (\lambda_i \mathbf{I} - \mathbf{A})^{+} \partial_{\theta} \mathbf{A}_{\theta} \mathbf{v}_i$
- 615 4.  $y = \log |\mathbf{A}|$ :  $\frac{\partial y}{\partial \theta} = \text{Tr}(\mathbf{A}^{-1} \partial_{\theta} \mathbf{A}_{\theta})$
- 616 5.  $\mathbf{y} = \text{Diag}(\mathbf{A})$ :  $\mathbf{w}^{\top} \frac{\partial \mathbf{y}}{\partial \theta} = \mathbf{w}^{\top} \text{Diag}(\partial_{\theta} \mathbf{A}_{\theta})$

617 In Figure 4 we show the practical benefits of our autograd rules. We take gradients of different  
 618 linear solves  $\mathbf{A}_{\theta}^{-1} \mathbf{b}$  that were derived using conjugate gradients (CG), where each solve required an  
 619 increasing number of CG iterations.

620 **C Algorithmic Details**

621 In this section we expand upon three different points introduced in the main paper. For the first point  
 622 we argue why SVRG leads to gradients with reduced variants. For the second points we display all  
 623 the iterative methods that we use as base algorithms in CoLA. Finally, for the third point we expand  
 624 upon CoLA’s strategy for dealing with the different numerical precisions that we support.

625 **C.1 SVRG**

626 In simplest form, SVRG [21] performs gradient descent with the varianced reduced gradient

$$\mathbf{w} \leftarrow \mathbf{w} - \eta(g_i(\mathbf{w}) - g_i(\mathbf{w}_0) + g(\mathbf{w}_0)) \tag{9}$$

627 where  $g_i$  represents the stochastic gradient evaluated at only a single element or minibatch of the sum,  
 628 and  $g(\mathbf{w}_0)$  is the full batch gradient evaluated at the anchor point  $\mathbf{w}_0$  which is recomputed at the end  
 629 of each epoch with an updated anchor.

630 With different loss functions, we can use this update rule to solve symmetric or asymmetric linear  
 631 systems, to compute the top eigenvectors or even find the nullspace of a matrix. Despite the fact that  
 632 the corresponding objectives are not strongly convex in the last two cases, it has been shown that  
 633 gradient descent and thus SVRG will converge at this exponential rate [51, 14]. Below we list the  
 gradients that enable us to solve different linear algebra problems: In each of the three cases listed

	Symmetric Solve $\mathbf{A}\mathbf{w} = \mathbf{b}$	Top- $k$ Eigenvectors $\mathbf{A}\mathbf{W} = \mathbf{W}\Lambda$	Nullspace $\mathbf{A}\mathbf{W} = 0$
$g_i(\mathbf{w})$	$\mathbf{A}_i\mathbf{w} - \mathbf{b}$	$-\mathbf{A}_i\mathbf{W} + \mathbf{W}\mathbf{W}^\top\mathbf{W}$ [51]	$\mathbf{A}_i\mathbf{W}$ [14]

Table 4: SVRG gradients for solving different linear algebra problems.

634 above, we can recognize that if the average of all the gradients  $g(w)$  is 0, then the corresponding  
 635 linear algebra solution has been recovered.  
 636

637 While it may seem that we need to take three complete passes through  $\{\mathbf{A}_i\}$  per SVRG epoch (due to  
 638 the three terms in Equation 9), we can reduce this cost to two complete passes exploiting the fact that  
 639 the gradients are linear in the matrix object, replacing  $\mathbf{A}_i\mathbf{W} - \mathbf{A}_i\mathbf{W}_0$  with  $\mathbf{A}_i(\mathbf{W} - \mathbf{W}_0)$  where  
 640 appropriate. In all of the Sum structure experiments where we leverage SVRG, the x-axis measures  
 641 the total number of passes through  $\{\mathbf{A}_i\}_{i=1}^m$ , two for each epoch for SVRG.

642 **C.2 Iterative methods**

643 In Table 5 we list the different iterative methods (base cases) that we use for different linear algebraic  
 644 operations as well as for different types of linear operators. As seen in Table 5, there are many  
 645 alternatives to our base cases, however we opted for algorithms that are known to be performant, that  
 646 are well-studied and that are popular amongst practitioners. A comprehensive explanation of our  
 647 bases cases and their alternatives can be found in Golub and Loan [19] and Saad [40].

648 **C.3 Lower precision linear algebra**

649 The accumulation of round-off error is usually the breaking point of several numerical linear algebra  
 650 (NLA) routines. As such, it is common to use precisions like float64 or higher, especially when  
 651 running these routines on a CPU. In contrast, in machine learning, lower precisions like float32 or  
 652 float16 are ubiquitously used because more parameters and data can be fitted into the GPU memory  
 653 (whose memory is usually much lower than CPUs) and because the MVMs can be done faster (the  
 654 CUDA kernels are optimized for operations on these precisions). Additionally, the round-off error  
 655 incurred on MVMs is not as detrimental when training machine learning models (as we are already  
 656 running noisy optimization algorithms) as when solving linear algebra problems (where round-off  
 657 error can lead us to poor solutions). Thus, it is an active area of research in NLA to derive routines  
 658 which utilize lower precisions than float64 or that mix precisions in order to achieve better runtimes  
 659 without a complete degradation of the quality of the solution.

660 In CoLA we take a two prong approach to deal with lower precisions in our NLA routines. First, we  
 661 incorporate additional variants of well-known algorithms that propagate less round-off error at the

Linear Algebra Op	Base Case	Alternatives
$\mathbf{Ax} = \mathbf{b}$ (asymmetric)	GMRES	BiCGSTAB, CR, QMR
$\mathbf{Ax} = \mathbf{b}$ (self-adjoint)	MINRES	GMRES
$\mathbf{Ax} = \mathbf{b}$ (PSD)	CG	GMRES
Eigs( $\mathbf{A}$ ) (asymmetric)	Arnoldi	IRAM, QR algorithm
Eigs( $\mathbf{A}$ ) (self-adjoint)	Lanczos	LOBPCG, Rayleigh-Ritz, Bi-Lanczos
$\mathbf{A}^+$	CG	CGS, LSQR, LGMRES
$\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^*$	Lanczos, rSVD	Jacobi-Davidson
$f(\mathbf{A})$ (self-adjoint)	SLQ	SVD, Rational Krylov Subspaces

Table 5: **CoLA’s base case iterative algorithm and some alternatives.** We now expand on the acronyms. GMRES: Generalized Minimum RESidual, BiCGSTAB: BiConjugate Gradient STABILized, CR: Conjugate Residuals, QMR: Quasi-Minimal Residual, MINRES: MINimum RESidual, CG: Conjugate Gradients, IRAM: Implicitly Restarted Arnoldi Method, LOBPCG: Locally Optimal Block Preconditioned Conjugate Gradients, Bi-Lanczos: Bidiagonal Lanczos, CGS: Conjugate Gradient Squared, LSQR: Least-Squares QR, LGMRES: Least-squares Generalized Minimum RESidual, SVD: Singular Value Decomposition, rSVD: randomized Singular Value Decomposition, and SLQ: Stochastic Lanczos Quadrature.

662 expense of requiring more computation, as seen in Figure 5. Second, we integrate novel variants of  
663 algorithms that are designed to be used on lower precisions such as the CG modification found in  
664 Maddox et al. [28]. We now discuss the first approach.

665 As discussed in Section C.2, there are two algorithms that are key for eigendecompositions. The first  
666 is Arnoldi (applicable to any operator), and the second is Lanczos (for symmetric operators) — where  
667 actually Lanczos can be viewed as a simplified version of Arnoldi. Central to these algorithms is the  
668 use of an orthogonalization step which is well-known to be a source of numerical instability. One  
669 approach to aggressively ameliorate the propagation of round-off error during orthogonalization is to  
670 use Householder projectors, which is the strategy that we use in CoLA. Given a unitary vector  $\mathbf{u}$ , a  
671 Householder projector (or Householder reflector) is defined as the following operator  $\mathbf{R} = \mathbf{I} - 2\mathbf{u}\mathbf{u}^*$ .  
672 When applied to a vector  $\mathbf{x}$  the result  $\mathbf{R}\mathbf{x}$  is basically a reflection of  $\mathbf{x}$  over the  $\mathbf{u}^\top$  space. To easily  
673 visualize this, suppose that  $\mathbf{x} \in \mathbb{R}^2$  and  $\mathbf{u} = \mathbf{e}_1$ . Hence,

$$\mathbf{R}\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} - 2 \begin{pmatrix} x_1 \\ 0 \end{pmatrix} = \begin{pmatrix} -x_1 \\ x_2 \end{pmatrix}$$

674 which is exactly the reflection of the vector across the axis generated by  $\mathbf{e}_2$ . Most notably,  $\mathbf{R}$  is unitary  
675  $\mathbf{R}\mathbf{R}^* = \mathbf{I}$  which can be easily verified from the definition. Being unitary is crucial as under the usual  
676 round-off error model, applying  $\mathbf{R}$  to another matrix  $\mathbf{A}$  does not worsen the already accumulated  
677 error  $\mathbf{E}$ . Mathematically,  $\|\mathbf{R}(\mathbf{A} + \mathbf{E}) - \mathbf{R}\mathbf{A}\| = \|\mathbf{R}\mathbf{E}\| = \|\mathbf{E}\|$ , where the last equality results from  
678 basic properties of unitary matrices. We are going to use Arnoldi as an example of how Householder  
679 projectors are used during orthogonalization. In Figure 5 we have an example of two different variants  
680 of Arnoldi present in CoLA. The implementations are notably different and also it is easy to see how  
681 Algorithm 2 is more expensive than Algorithm 1. First, note that for Algorithm 2 we have two for  
682 loops (line 6 and line 8) whereas for Algorithm 1 we only have one (line 4-6). Worse, the two for  
683 loops in Algorithm 2 require more flops than the only for loop in Algorithm 1. Note that we do not  
684 always favor the more expensive but robust implementation of an algorithm as in some cases, like  
685 when running GMRES, the round-off error is not as impactful to the quality of the solution, and  
686 shorter runtimes are actually more desirable.

## 687 D Experimental Details

688 In this section we expand upon the details of all the experiments ran in the paper. Such details include  
689 the datasets that were used, the hyperparameters of different algorithms and the specific choices

---

**Algorithm 1** Arnoldi iteration

---

```
1: Inputs:  $\mathbf{A}$ ,  $\mathbf{q}_0 = \boldsymbol{\nu}_0 / \|\boldsymbol{\nu}_0\|$  where possibly  $\boldsymbol{\nu}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , maximum number of iterations  $T$  and tolerance  $\epsilon \in (0, 1)$ .
2: for  $j = 0$  to  $T - 1$  do
3:    $\boldsymbol{\nu}_{j+1} \leftarrow \mathbf{A}\mathbf{q}_j$ 
4:   for  $i = 0$  to  $j$  do
5:      $h_{i,j} = \mathbf{q}_i^* (\mathbf{A}\mathbf{q}_j)$ 
6:      $\boldsymbol{\nu}_{j+1} \leftarrow \boldsymbol{\nu}_{j+1} - h_{i,j} \mathbf{q}_i$ 
7:   end for
8:    $h_{j+1,j} = \|\boldsymbol{\nu}_{j+1}\|$ 
9:   if  $h_{j+1,j} < \epsilon$  then
10:    stop
11:  else
12:     $\mathbf{q}_{j+1} = \boldsymbol{\nu}_{j+1} / h_{j+1,j}$ 
13:  end if
14: end for
15: return  $\mathbf{H}, \mathbf{Q} = (\mathbf{q}_0 | \dots | \mathbf{q}_{T-1} | \mathbf{q}_T)$ 
```

---

---

**Algorithm 2** Householder Arnoldi iteration

---

```
1: Inputs:  $\mathbf{A}$ ,  $\boldsymbol{\nu}_0 \neq \mathbf{0}$  where possibly  $\boldsymbol{\nu}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , and maximum number of iterations  $T$ .
2: for  $j = 0$  to  $T$  do
3:    $\mathbf{u}_j = \text{GET\_HOUSEHOLDER\_VEC}(\boldsymbol{\nu}_j, j)$ 
4:    $\mathbf{R}_j = \mathbf{I} - 2\mathbf{u}_j \mathbf{u}_j^*$ 
5:    $\mathbf{h}_j = \mathbf{R}_j \boldsymbol{\nu}_j$ 
6:    $\mathbf{q}_j = \mathbf{R}_0 \cdots \mathbf{R}_j \mathbf{e}_{j+1}$ 
7:   if  $j < T$  then
8:      $\boldsymbol{\nu}_{j+1} = \mathbf{R}_j \cdots \mathbf{R}_0 (\mathbf{A}\mathbf{q}_j)$ 
9:   end if
10: end for
11: return  $\mathbf{H}, \mathbf{Q} = (\mathbf{q}_0 | \dots | \mathbf{q}_T)$ 
12: function GET_HOUSEHOLDER_VEC( $\mathbf{w}, k$ )
13:    $u_i = 0$  for  $i < k$  and  $u_i = w_i$  for  $i > k$ .
14:    $u_k = w_k - \|\mathbf{w}\|$ 
15:   return  $\mathbf{u}$ 
16: end function
```

---

Figure 5: Different versions of the same algorithm, but the Householder variant being more numerically robust.

690 of algorithms used both for CoLA but also for the alternatives. We run each of the experiments 3  
691 times and compute the mean dropping the first observation (as usually the first run contains some  
692 compiling time much is not too large). We do not display the standard deviation as those numbers  
693 are imperceptible for each experiment. In terms of hardware, the CPU experiments were run on an  
694 Intel(R) Core(TM) i5-9600K CPU @ 3.70GHz and the GPU experiments were run on a NVIDIA  
695 GeForce RTX 2080 Ti.

## 696 D.1 Datasets

697 Below we enumerate the datasets that we used in the various applications. Most of the datasets are  
698 sourced from the University of California at Irvine’s (UCI) Machine Learning Respository that can  
699 be found here: <https://archive.ics.uci.edu/ml/datasets.php>. Also, a community repo  
700 hosting these UCI benchmarks can be found here: [https://github.com/treforevans/uci\\_](https://github.com/treforevans/uci_datasets)  
701 [datasets](https://github.com/treforevans/uci_datasets) (we have no affiliation).

- 702 1. *Elevators*. This dataset is a modified version of the *Ailerons* dataset, where the goal is to  
703 to predict the control action on the ailerons of the aircraft. This UCI dataset consists of  
704  $N = 14\text{K}$  observations and has  $D = 18$  dimensions.
- 705 2. *Kin40K*. The full name of this UCI dataset is *Statlog (Shuttle) Data Set*. This dataset contains  
706 information about NASA shuttle flights and we used a subset that consists of  $N = 40\text{K}$   
707 observations and has  $D = 8$  dimensions.
- 708 3. *Buzz*. The full name of this UCI dataset is *Buzz in social media*. This dataset consists of  
709 examples of buzz events from Twitter and Tom’s Hardware. We used a subset consisting of  
710  $N = 430\text{K}$  observations and has  $D = 77$  dimensions.
- 711 4. *Song*. The full name of this UCI dataset is *YearPredictionMSD*. This dataset consists of  
712  $N = 386.5\text{K}$  observations and it has  $D = 90$  audio features such as 12 timbre average  
713 features and 78 timbre covariance features.
- 714 5. *cit-HepPh*. This dataset is based on arXiv’s HEP-PH (high energy physics phenomenology)  
715 citation graph and can be found here: [https://snap.stanford.edu/data/cit-HepPh.](https://snap.stanford.edu/data/cit-HepPh.html)  
716 [html](https://snap.stanford.edu/data/cit-HepPh.html). The dataset covers all the citations from January 1993 to April 2003 of  $|V| = 34, 549$   
717 papers, ultimately containing  $|E| = 421, 578$  directed edges. The notion of relationship  
718 that we used in our spectral clustering experiment creates a connection between two papers

719 when at least one cites another (undirected symmetric graph). Therefore the dataset that we  
720 used has the same number of nodes but instead  $|E| = 841,798$  undirected edges.

## 721 D.2 Compositional experiments

722 This section pertains to the experiments of Section 3.2 displayed in Figure 1. We now elaborate on  
723 each of Figure 1’s panels.

- 724 (a) The multi-task GP problem exploits the structure of the following Kronecker operator  
725  $\mathbf{K}_T \otimes \mathbf{K}_X$ , where  $\mathbf{K}_T$  is a kernel matrix containing the correlation between the tasks and  
726  $\mathbf{K}_X$  is a RBF kernel on the data. For this experiment, we used a synthetic Gaussian dataset  
727 where the train data  $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_D)$  which has dimension  $D = 33$ ,  $N = 1\text{K}$  and we used  
728  $T = 11$  tasks (where the tasks basically set the size of  $\mathbf{K}_T$ ). We used conjugate gradients  
729 (CG) as the iterative method, where we set the hyperparameters to a tolerance of  $10^{-6}$  and  
730 to a maximum number of iterations to 1K. We used the exact same hyperparameters for  
731 CoLA.
- 732 (b) For the bi-poisson problem we set up the maximum grid to be  $N = 1000^2$ . Since this PDE  
733 problem involves solving a symmetric linear system, we used CG as the iterative method  
734 with a tolerance of  $10^{-11}$  and a maximum number of iterations of 10K. The previous  
735 parameters also apply for CoLA. We note that PDE problems are usually solved to higher  
736 tolerances as the numerical error compounds as we advance the PDE.
- 737 (c) For the EMLP experiment we consider solving the equivariance constraints to find the  
738 equivariant linear layers of a graph neural network with 5 nodes. To solve this problem,  
739 we need to find the nullspace of a large structured constraint matrix. We use the uniformly  
740 channel heuristic from [14] which distributes the  $N$  channels across tensors of different  
741 orders. We consider our approach which exploits the block diagonal structure, separating  
742 the nullspaces into blocks, as opposed to the direct iterative approach exploiting only the  
743 fast MVMs of the constraint matrix. We use a tolerance of  $10^{-5}$ .

## 744 D.3 Sum structure experiments

745 This section pertains to the experiments of Section 3.3 contained in Figure 2. We now elaborate on  
746 each of Figure 2’s panels.

- 747 (a) In this experiment we computed the first principal component of the *Buzz* dataset. For the  
748 iterative method we used power iteration with a maximum number of iterations of 300 and a  
749 stop tolerance of  $10^{-7}$ . CoLA used SVRG also with the same stop tolerance and maximum  
750 number of iterations. Additionally, we set SVRG’s batch size to 10K and the learning rate  
751 to 0.0008. We note that a single power iteration roughly contains  $43/2 = 21.5$  times more  
752 MVMs than a single iteration of SVRG. In this particular case, the length of the sum is given  
753 by the number of observations and therefore SVRG uses  $430/10 = 43$  times less elements  
754 per iteration, where 10 comes from the 10K batch size. Finally, the 2 is explained by noting  
755 that SVRG incurs in a full sum update on every epoch.
- 756 (b) In this experiment we trained a GP by estimating the covariance RBF kernel with  $J = 1\text{K}$   
757 random Fourier features (RFFs). The hyperparameters for the RBF kernel are the following:  
758 length scale ( $\ell = 0.1$ ), output scale ( $a = 1$ ) and likelihood noise ( $\sigma^2 = 0.1$ ). Moreover, we  
759 used CG as the iterative solver with a tolerance of  $10^{-8}$  and 100 as the maximum number  
760 of iterations (the convergence took much less iterations than the max). For SVRG we used  
761 the same tolerance but set the maximum number of iterations to 10K, a batch size of 100  
762 and learning rate of 0.004. We note that a single CG iteration roughly contains  $10/2 = 5$   
763 times more MVMs than a single iteration of SVRG. In this particular case, the length of the  
764 sum is given by the number of RFFs and therefore SVRG uses  $1000/100 = 10$  times less  
765 elements per iteration, where 100 comes from the batch size.
- 766 (c) In this experiment we implemented the Neural-IVP method from Finzi et al. [15]. We  
767 consider the time evolution of a wave equation in two spatial dimensions. At each integrator  
768 step, a linear system  $\mathbf{M}(\theta)\dot{\theta} = F(\theta)$  must be solved to find  $\dot{\theta}$ , for a  $d = 12\text{K} \times 12\text{K}$   
769 dimensional matrix. While Finzi et al. [15] use conjugate gradients to solve the linear  
770 system, we demonstrate the advantages of using SVRG, as  $\mathbf{M}(\theta) = \frac{1}{m} \sum_{i=1}^m M_i(\theta)$  is a

771 sum over the evaluation at  $m = 50\text{K}$  distinct sample locations within the domain. In this  
772 experiment we use a batch size of 500 for SVRG, and employ rank 250 randomized Nyström  
773 preconditioning for both SVRG and the iterative CG baseline.

#### 774 D.4 Applications

775 This section pertains to the experiments of Section 4 displayed in Figure 3. We now elaborate on  
776 each of Figure 3’s panels.

- 777 (a) In this experiment we compute 5, 10 and 20 PCA components for the *Buzz* dataset. We  
778 compared against `sklearn` which uses the Lanczos algorithm through the fast Fortran-based  
779 ARPACK numerical library. In this case, CoLA uses randomized SVD [31] with a rank 3000  
780 approximation.
- 781 (b) In this experiment we fit a Ridge regression on the *Song* dataset with a regularization  
782 coefficient set to 0.1. We compared against `sklearn` using their fastest least-square solver  
783 `lsqr` with a tolerance of  $10^{-4}$ . In this case, CoLA uses CG with the same tolerance and  
784 with a maximum number of iterations set to 1K. Additionally, we ran CoLA using CPU and  
785 GPU whereas we used only CPU for `sklearn` as it has no GPU support. We observe how  
786 in the arguably most popular ML method, CoLA is able to beat a leading package such as  
787 `sklearn`.
- 788 (c) In this experiment we fit a GP with a RBF kernel on two datasets: *Elevators* and *Kin40K*.  
789 We only used up to 20K observations from *Kin40K* as that was the maximum number of  
790 observations that would fit the GPU memory without needing to partition the MVMs. We  
791 compare against `GPyTorch` which uses CG and stochastic Lanczos quadrature (SLQ) to  
792 compute and optimize the negative log-marginal likelihood (loss function). Both experiments  
793 were run on a GPU for 100 iterations using Adam as an optimizer with learning rate of  
794 0.1 with the default values of  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$ . Additionally, for both `GPyTorch`  
795 and CoLA, the CG tolerance was set to  $10^{-4}$  with a maximum number of CG iterations of  
796 250 and 20 probes were used for SLQ. Note that both CoLA and `GPyTorch` have similar  
797 throughputs, for example `GPyTorch` runs a 100 iterations on *Elevators* on 43 seconds  
798 whereas CoLA runs a 100 iterations on 49 seconds. When training a GP, we solve a block of  
799 11 linear systems (1 based on  $y$  and 10 based on random probes) where one key difference  
800 is that the CG solver for `GPyTorch` has a stopping criteria based on the convergence of  
801 the mean solves whereas CoLA has a stopping criteria based on the convergence of all the  
802 solves.
- 803 (d) In this experiment we run spectral clustering on the *cit-HepPh* dataset using an embedding  
804 size of 8 and also 8 clusters for k-means (with only 1 run of k-means after estimating the  
805 embeddings). We compare against `sklearn` using two different solvers, one based on  
806 Lanczos iterations using ARPACK and another using an Algebraic Multi-Graph solver AMG. In  
807 this case, CoLA also uses Lanczos iterations with a default tolerance of  $10^{-6}$ . We see how  
808 `sklearn`’s AMG solver runs faster than CoLA’s but this is mostly the algorithmic constants  
809 as they have similar asymptotical behavior (similar slopes).
- 810 (e) In this experiment we solve the Schrödinger equation to find the energy levels of the  
811 hydrogen atom on a 3-dimensional finite difference grid with up to  $N = 5\text{K}$  points. In  
812 order to handle the infinite spatial extent, we compactify the domain by applying the arctan  
813 function. Under this change of coordinates, the Laplacian has a different form, and hence  
814 the matrix forming the discretized Hamiltonian is no longer symmetric. We compare against  
815 `SciPy`’s Arnoldi implementation with 20 iterations where CoLA also uses Arnoldi with the  
816 same number of iterations. Surprisingly, CoLA’s JAX jitted code has a competitive runtime  
817 when compare to `SciPy`’s runtime using ARPACK.
- 818 (f) In this experiment we solve a minimal surface problem on a grid of maximum size of  
819  $N = 100^2$  points. To solve this problem we have to run `NetworkX`’s `Rhaphson` where each  
820 inner step involves a linear solve of an asymmetric operator. We compare against `SciPy`’s  
821 `GMRES` implementation as well as JAX’s integrated version of `SciPy`. The main difference  
822 between the two is that `SciPy` calls the fast and highly-optimized ARPACK library whereas  
823 `SciPy` (JAX) has its only Python implementation of `GMRES` which only uses JAX’s  
824 primitives (equally as it is done in CoLA). The tolerance for this experiment was  $5\text{e-}3$ .



```
887 def inverse(A: Diagonal, **kwargs) -> Diagonal:
888     return Diagonal(1. / A.diag)
889
890
891 @dispatch
892 def inverse(A: Unitary, **kwargs) -> Unitary:
893     return Unitary(A.H)
```