## Practical and Efficient Time Integration and Kronecker Product Solvers

Jed Brown jed.brown@colorado.edu (CU Boulder and ANL) Collaborators: Debojyoti Ghosh (LLNL), Matt Normile (CU), Martin Schreiber (Exeter)

Preconditioning, 2017-08-01 This talk: https: //jedbrown.org/files/20170801-FastKronecker.pdf

< □ > < 同 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <



Floating Point Operations per Byte, Double Precision

[c/o Karl Rupp]

### 2017 HPGMG performance spectra



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへで

### **Motivation**

#### Hardware trends

- Memory bandwidth a precious commodity (8+ flops/byte)
- Vectorization necessary for floating point performance
- Conflicting demands of cache reuse and vectorization
- Can deliver bandwidth, but latency is hard
- Assembled sparse linear algebra is doomed!
  - Limited by memory bandwidth (1 flop/6 bytes)
  - No vectorization without blocking, return of ELLPACK
- Spatial-domain vectorization is intrusive
  - Must be unassembled to avoid bandwidth bottleneck

- Whether it is "hard" depends on discretization
- Geometry, boundary conditions, and adaptivity

Sparse linear algebra is dead (long live sparse ...)

- Arithmetic intensity < 1/4</li>
- Idea: multiple right hand sides

 $\frac{(2k \text{ flops})(\text{bandwidth})}{\text{sizeof}(\text{Scalar}) + \text{sizeof}(\text{Int})}, \quad k \ll \text{avg. nz/row}$ 

- Problem: popular algorithms have nested data dependencies
  - Time step
     Nonlinear solve
     Krylov solve
     Preconditioner/sparse matrix
- Cannot parallelize/vectorize these nested loops
- Can we create new algorithms to reorder/fuse loops?
  - Reduce latency-sensitivity for communication
  - Reduce memory bandwidth (reuse matrix while in cache)

Sparse linear algebra is dead (long live sparse ...)

- Arithmetic intensity < 1/4</li>
- Idea: multiple right hand sides

 $\frac{(2k \text{ flops})(\text{bandwidth})}{\text{sizeof}(\text{Scalar}) + \text{sizeof}(\text{Int})}, \quad k \ll \text{avg. nz/row}$ 

- Problem: popular algorithms have nested data dependencies
  - Time step
     Nonlinear solve
     Krylov solve
     Preconditioner/sparse matrix
- Cannot parallelize/vectorize these nested loops
- Can we create new algorithms to reorder/fuse loops?
  - Reduce latency-sensitivity for communication
  - Reduce memory bandwidth (reuse matrix while in cache)

### Attempt: s-step methods in 3D



- Limited choice of preconditioners (none optimal, surface/volume)
- Amortizing message latency is most important for strong-scaling
- s-step methods have high overhead for small subdomains

### Attempt: distribute in time (multilevel SDC/Parareal)



- PFASST algorithm (Emmett and Minion, 2012)
- Zero-latency messages (cf. performance model of s-step)
- Spectral Deferred Correction: iterative, converges to IRK (Gauss, Radau, ...)
- Stiff problems use implicit basic integrator (synchronizing on spatial communicator)

### Problems with SDC and time-parallel



c/o Matthew Emmett, parallel compared to sequential SDC

- Iteration count not uniform in s; efficiency starts low
- Low arithmetic intensity; tight error tolerance (cf. Crank-Nicolson)
- Parabolic space-time also works, but comparison flawed

### Runge-Kutta methods



- General framework for one-step methods
- Diagonally implicit: A lower triangular, stage order 1 (or 2 with explicit first stage)
- Singly diagonally implicit: all A<sub>ii</sub> equal, reuse solver setup, stage order 1
- ► If A is a general full matrix, all stages are coupled, "implicit RK"

### Implicit Runge-Kutta



- Excellent accuracy and stability properties
- Gauss methods with s stages
  - ▶ order 2*s*, (*s*, *s*) Padé approximation to the exponential

◆□▶ ◆帰▶ ◆ヨ▶ ◆ヨ▶ = ● ののの

- A-stable, symplectic
- Radau (IIA) methods with s stages
  - ▶ order 2s 1, A-stable, L-stable
- Lobatto (IIIC) methods with s stages
  - order 2s 2, A-stable, L-stable, self-adjoint
- Stage order s or s+1

### Method of Butcher (1976) and Bickart (1977)

Newton linearize Runge-Kutta system at u\*

$$Y = u^n + hAF(Y)$$
  $[I_s \otimes I_n + hA \otimes J(u^*)]\delta Y = RHS$ 

Solve linear system with tensor product operator

$$\hat{G} = S \otimes I_n + I_s \otimes J$$

where  $S = (hA)^{-1}$  is  $s \times s$  dense,  $J = -\partial F(u)/\partial u$  sparse

- SDC (2000) is Gauss-Seidel with low-order corrector
- Butcher/Bickart method: diagonalize  $S = V\Lambda V^{-1}$ 
  - $\Lambda \otimes I_n + I_s \otimes J$
  - s decoupled solves
  - Complex eigenvalues (overhead for real problem)

## III conditioning

$$A = V \Lambda V^{-1}$$

#### Conditioning of eigenbasis: gauss



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ ─臣 ─のへで

### Skip the diagonalization



- Accessing memory for J dominates cost
- Irregular vector access in application of J limits vectorization
- Permute Kronecker product to reuse J and make fine-grained structure regular
- Stages coupled via register transpose at spatial-point granularity
- Same convergence properties as Butcher/Bickart

### MatKAIJ: "sparse" Kronecker product matrices

$$G=I_n\otimes S+J\otimes T$$

- ► J is parallel and sparse, S and T are small and dense
- More general than multiple RHS (multivectors)
- Compare  $J \otimes I_s$  to multiple right hand sides in row-major
- Runge-Kutta systems have T = I<sub>s</sub> (permuted from Butcher method)
- Stream J through cache once, same efficiency as multiple RHS
- Unintrusive compared to spatial-domain vectorization or s-step

### Convergence with point-block Jacobi preconditioning

#### 3D centered-difference diffusion problem

Method	order	nsteps	Krylov its.	(Average)
Gauss 1	2	16	130	(8.1)
Gauss 2	4	8	122	(15.2)
Gauss 4	8	4	100	(25)
Gauss 8	16	2	78	(39)

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ ▲ □ ● ● ● ●

### We really want multigrid

- Coarse operator:  $I_n \otimes S + (RJP) \otimes I_s$
- Larger time steps
- GMRES(2)/point-block Jacobi smoothing
- FGMRES outer

Method	order	nsteps	Krylov its.	(Average)
Gauss 1	2	16	82	(5.1)
Gauss 2	4	8	64	(8)
Gauss 4	8	4	44	(11)
Gauss 8	16	2	42	(21)

▲ロト ▲ □ ト ▲ □ ト ▲ □ ト ● ● の Q ()

### Toward a better AMG for IRK/tensor-product systems



Start with 
$$\hat{R} = R \otimes I_s$$
,  $\hat{P} = P \otimes I_s$ 

$$G_{\text{coarse}} = \hat{R}(I_n \otimes S + J \otimes I_s)\hat{P}$$

- Imaginary component slows convergence
- Can we use a Kronecker product interpolation?
- Rotation on coarse grids (connections to shifted Laplacian)

イロト 人間 とくほ とくほ とう

-

### Why implicit is silly for waves

- Implicit methods require an implicit solve in each stage.
- Time step size proportional to CFL for accuracy reasons.
- Methods higher than first order are not unconditionally strong stability preserving (SSP; Spijker 1983).
  - Empirically, c<sub>eff</sub> ≤ 2, Ketcheson, Macdonald, Gottlieb (2008) and others

< □ > < 同 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

- Downwind methods offer to bypass, but so far not practical
- Time step size chosen for stability
  - Increase order if more accuracy needed
  - Large errors from spatial discretization, modest accuracy
- My goal: need less memory motion per stage
  - Better accuracy, symplecticity nice bonus only
  - Cannot sell method without efficiency

### Implicit Runge-Kutta for advection

Table: Total number of iterations (communications or accesses of *J*) to solve linear advection to t = 1 on a 1024-point grid using point-block Jacobi preconditioning of implicit Runge-Kutta matrix. The relative algebraic solver tolerance is  $10^{-8}$ .

Family	Stages	Order	Iterations
Crank-Nicolson/Gauss	1	2	3627
Gauss	2	4	2560
Gauss	4	8	1735
Gauss	8	16	1442

< □ > < 同 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

- Naive centered-difference discretization
- Leapfrog requires 1024 iterations at CFL=1
- This is A-stable (can handle dissipation)

### **Diagonalization revisited**

$$(I \otimes I - hA \otimes L)Y = (\mathbf{1} \otimes I)u_n \tag{1}$$

$$u_{n+1} = u_n + h(b^T \otimes L)Y$$
<sup>(2)</sup>

Sac

• eigendecomposition 
$$A = V \Lambda V^{-1}$$
  
 $(V \otimes I)(I \otimes I - h \Lambda \otimes L)(V^{-1} \otimes I)Y = (\mathbf{1} \otimes I)u_n.$ 

- Find diagonal W such that  $W^{-1}\mathbf{1} = V^{-1}\mathbf{1}$
- Commute diagonal matrices

$$(I \otimes I - h\Lambda \otimes L) \underbrace{(WV^{-1} \otimes I)Y}_{Z} = (\mathbf{1} \otimes I)u_n.$$

• Using  $\tilde{b}^T = b^T V W^{-1}$ , we have the completion formula

$$u_{n+1} = u_n + h(\tilde{b}^T \otimes L)Z.$$

- $\Lambda, \tilde{b}$  is new diagonal Butcher table
- Compute coefficients offline using extended precision to handle

### **Exploiting realness**

Eigenvalues come in conjugate pairs

$$A = V\Lambda V^{-1}$$

For each conjugate pair, create unitary transformation

$$T = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1\\ i & -i \end{bmatrix}$$

► Real 2 × 2 block diagonal *D*; real  $\tilde{V}$  (with appropriate phase)

$$A = (VT^*)(T\Lambda T^*)(TV^{-1}) = \tilde{V}\tilde{D}V^{-1}$$

- Yields new block-diagonal Butcher table D, b.
- Halve number of stages using identity

$$\overline{(\alpha+J)^{-1}u}=(\overline{\alpha}+J)^{-1}u$$

Solve one complex problem per conjugate pair, then take twice the real part.

### **REXI:** Rational approximation of exponential

$$u(t)=e^{Lt}u(0)$$

Haut, Babb, Martinsson, Wingate; Schreiber and Loft

$$(\alpha \otimes I + hI \otimes L)Y = (\mathscr{W} \otimes I)u_n$$
$$u_{n+1} = (\beta^T \otimes I)Y.$$

- $\alpha$  is complex-valued diagonal,  $\beta$  is complex
- Constructs rational approximations of Gaussian basis functions, target (real part of) e<sup>it</sup>
- REXI is a Runge-Kutta method: can convert via "modified Shu-Osher form"
  - Developed for SSP (strong stability preserving) methods
  - Ferracina, Spijker (2005), Higueras (2005)
  - Yields diagonal Butcher table  $A = -\alpha^{-1}, b = -\alpha^{-2}\beta$

### Abscissa for RK and REXI methods



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへで

### Stability regions





◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへで

### Computational Performance (SWE on the plane)

# Spectral solver with RK4 time stepping method vs. REXI with spectral solver



Resolution  $= 128^2$ 

More than one order of magnitude **faster** with **similar accuracy** 

Proof of concept that REXI works with spectral methods

イロト 不得 とうほう イヨン

Computed on Linux Cluster, LRZ / Technical University of Munich M. Schreiber, P. Peixoto, TS. Haut, BA. Wingate - Beyond spatial scalability limitations with a massively parallel method for linear oscillatory problems, in International Journal of High Performance Computing Applications

18/63

### Outlook on Kronecker product solvers

 $I \otimes S + T \otimes J$ 

- (Block) diagonal S is usually sufficient
- Best opportunity for "time parallel" (for linear problems)
  - Is it possible to beat explicit wave propagation with high efficiency?
- Same structure for stochastic Galerkin and other UQ methods
- IRK unintrusively offers bandwidth reuse and vectorization
- Need polynomial smoothers for IRK spectra
- Change number of stages on spatially-coarse grids (p-MG, or even increase)?
- Experiment with SOR-type smoothers
  - Prefer point-block Jacobi in smoothers for spatial parallelism
- Possible IRK correction for IMEX (non-smooth explicit function)
- PETSc implementation (works in parallel, hardening in progress)
- Thanks to DOE ASCR