

# RSLAB

## A Pure-Rust Sparse Direct Solver: Algorithms, A-Priori Predictors, and a Learned Configuration Tuner

Milan Rother

July 2, 2026

### Abstract

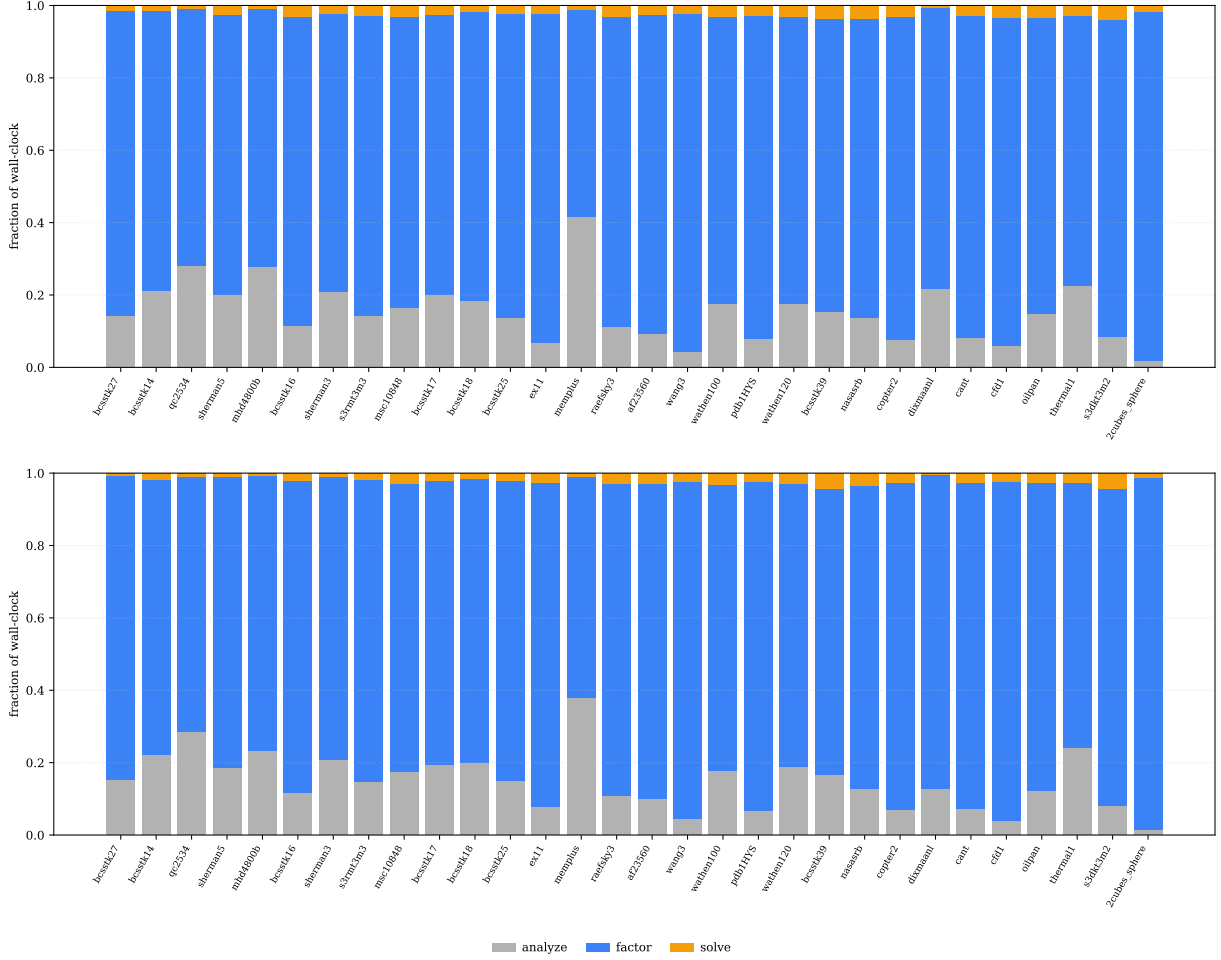
RSLAB is a pure-Rust sparse direct solver. It provides a supernodal complex-symmetric  $\mathbf{LDL}^\top$  factorization with Bunch–Kaufman pivoting and an unsymmetric multifrontal LU, is generic over the scalar type (`f64`, `f32`, `Complex<f64>`, `Complex<f32>`), and links no external numeric library (no BLAS, LAPACK, or METIS): the SIMD dense kernels, the fill-reducing orderings, and the factorization are all Rust. The analyze phase computes, in addition to the symbolic factorization, a set of a-priori predictors that are exact functions of the matrix structure: a per-path peak-memory bound (separate for the left-looking and the multifrontal schedule), a geometric-flop work estimate, a runtime estimate from a hardware calibration, and a thread-count recommendation. The solver exposes its analyze, factor, and kernel parameters through one flat settings interface, and a learned configuration tuner selects a setting from the matrix structure, constrained by a deterministic guard stack that keeps its predicted memory below the untuned default and falls back to the default outside the training range. The numeric factorization is bit-identical regardless of the worker count, and the resource planner is a pure function of its inputs, which makes batched execution reproducible and its memory and runtime cost computable in advance. This report documents the algorithms, the predictors, the tunable stages, the determinism properties, and the solver-in-the-loop execution model, and reports the solver against `faer`, MKL PARDISO, and SuperLU over the SuiteSparse collection.

## 1 Scope

RSLAB factorizes and solves  $Ax = b$  for sparse  $A$  that is either complex-symmetric (routed to  $\mathbf{LDL}^\top$ ) or general unsymmetric (routed to LU). The primary target is the complex-symmetric systems of frequency-domain electromagnetic and FEM analysis, including curl–curl (Nédélec edge-element) Maxwell problems with lossy media or PML, which are complex-symmetric rather than Hermitian. The implementation is stable Rust with no C/Fortran dependency and no BLAS/LAPACK/METIS linkage, and is exposed to Python through a PyO3 layer.

## 2 Architecture

The solver follows the standard three-phase sparse-direct flow. *Analyze* orders the matrix, builds the elimination tree, detects and amalgamates supernodes, computes column counts, and produces the a-priori predictors (Sec. 6) and the structural feature vector (Sec. 6). *Factor* performs the numeric  $\mathbf{LDL}^\top$  or LU on the ordered, equilibrated matrix. *Solve* runs the triangular substitutions. Analyze is separated from factor at the API level: a symbolic object (`LdlSymbolic` or `LuSymbolic`) is produced once and can factor many matrices that share its sparsity pattern, which is the basis of the solver-in-the-loop model (Sec. 10). Across the corpus the numeric factor dominates the wall-clock, with analyze and solve small by comparison (Fig. 1), so the tuning and the parallelism target the factor.



**Figure 1:** Analyze / factor / solve wall-clock split per matrix, for the left-looking and multifrontal paths. The numeric factor dominates.

The code is a Cargo workspace (Table 1). The ordering crates share a common quotient-graph elimination engine; the numeric core holds the two factorization paths and the dense SIMD kernels; the diagnostics, analysis, tuning, and `auto_tune` modules hold the predictors, the feature extraction, the resource planner, and the learned tuner respectively.

### 3 Fill-reducing orderings

The ordering determines fill and therefore both factor time and memory. RSLAB implements three ordering families behind an `OrderingMethod` enum and an adaptive dispatcher.

#### 3.1 Approximate minimum degree

AMD [1] is a quotient-graph elimination that selects the pivot of least approximate external degree. The workspace holds the element and variable lists in a single sliding index array `iw`, with element pointers `pe`, list lengths `len/elen`, supervariable counts `nv`, running degrees `degree`, a generation-counter mark array `w`, and degree-indexed LIFO bucket lists `head/next/last`; when `iw` fills, an in-place garbage collection compacts the live lists. The elimination loop repeatedly selects the minimum-degree pivot, forms a new element by merging its element list with its variable tail (in place when the pivot has no prior elements, out of place otherwise), and updates the affected variables' approximate degrees. During the update it performs aggressive element absorption, mass elimination of variables that become interchangeable with the pivot, and supervariable detection

Crate / module	Responsibility
<code>rslab-ordering-core</code>	Shared quotient-graph elimination engine (workspace, sliding index allocator with garbage collection, supervariable detection)
<code>rslab-amd</code>	Approximate minimum degree
<code>rslab-amf</code>	Approximate minimum fill (HAMF4)
<code>rslab-metis</code>	Multilevel nested dissection (iterative driver, König separators)
<code>symbolic</code>	Elimination tree, supernode amalgamation, column counts, ordering dispatch
<code>numeric</code>	Supernodal $\mathbf{LDL}^T$ , multifrontal LU, panel-freeing, dense SIMD kernels, scoped pools
<code>diagnostics</code>	Per-path a-priori memory and flop estimates
<code>analysis</code>	Structural feature extraction, thread-count predictor
<code>tuning</code>	Hardware calibration, runtime estimate, budget-aware resource planner
<code>auto_tune</code>	Embedded performance model and guard stack (pure-Rust inference)
<code>python</code>	PyO3 bindings

**Table 1:** Workspace layout.

by hashing variables with identical adjacency and merging the indistinguishable ones. Variables of degree exceeding  $\max(16, \lfloor \alpha \sqrt{n} \rfloor)$  with  $\alpha = 10$  are deferred to the tail as dense rows. AMD is the default for very large, very sparse patterns, where the separator overhead of nested dissection does not pay off.

### 3.2 Approximate minimum fill

AMF [2, 3] replaces the degree score with an approximate *fill* score. Each supervariable  $i$  carries

$$\text{RMF}(i) = \frac{\deg(i) (\deg(i) - 1 + 2 \deg_{me}) - \text{WF}(i)}{nv(i) + 1}, \quad (1)$$

where  $\text{WF}(i)$  is the accumulated working fill, updated per eliminated element by a surface term  $\text{dext} (2 \deg - \text{dext} - 1)$  and combined across supervariables as  $\text{wf}_4 + 2 nv_i \text{wf}_3$ . The working fill is accumulated in 64-bit integers because the product is  $\mathcal{O}(n^2)$  and overflows 32-bit for  $n \gtrsim 46\,000$ . Scores are placed in a bucket array of length  $2n + 2$  with exact fine buckets  $[0, n]$  and coarse buckets of stride  $\max(n/8, 1)$  above; the pivot is the minimum exact score. Supervariable merges take the larger of the two working-fill values. AMF is the HAMF4 variant and runs aggressive absorption unconditionally; it is the default for  $n \leq 10\,000$ , where it stays within  $1.10\times$  of the reference HAMF4 fill on a 183 277-matrix oracle suite.

### 3.3 Nested dissection

The nested-dissection ordering [4, 5] is a multilevel recursive bisection driven by an explicit work stack of (subgraph, vertex map, offset) items rather than recursion, so deep dissection trees cannot overflow the call stack. Each subgraph is bisected by coarsening the graph, computing an initial bisection on the coarsest level, refining it with Fiduccia–Mattheyses [6], and uncoarsening with projection and refinement at each level. The vertex separator is extracted from the edge-cut bipartition by a minimum vertex cover, computed by König’s theorem: a maximum matching on the bipartite crossing-edge graph, then the alternating-path reachable set gives the cover. Two guards bound the recursion: subgraphs below `nd_to_amd_switch` = 200 vertices are ordered directly by AMD, and a split where one side holds  $\geq 90\%$  of the vertices falls back to an AMD leaf rather than recursing on a graph without a good separator. Nested dissection is the default for  $n > 10\,000$ .

### 3.4 Dispatch

`OrderingMethod::Auto` routes by cheap structural predicates. A very large and sparse pattern ( $n > 100\,000$ , average degree  $< 5$ ) goes to AMD. An arrow or bordered pattern goes to AMF,

which defers the dense border to the trailing block where nested dissection would smear it across separators; the detector flags a pattern where at least 20% of the nonzeros are concentrated in fewer than 5% of the columns above a degree floor of  $\max(64, 8\overline{\text{deg}})$ . Otherwise the size rule above applies. `OrderingMethod::AutoRace` instead runs the full symbolic factorization for each concrete candidate and keeps the one with the smallest factor nnz, at about  $4\times$  the single-pass analysis cost.

## 4 Symbolic analysis

From the ordered pattern RSLAB builds the elimination tree and detects fundamental supernodes, which are maximal runs of columns whose row structure differs only by the eliminated column. Small supernodes are amalgamated to widen the dense fronts that feed the BLAS-3 kernels. Nodes with fewer than `nemin` = 16 columns [7] are merged with their parent; the merge order uses a column renumbering that re-postorders the tree so a desired parent-child merge becomes structurally adjacent, which is what makes amalgamation effective on bushy assembly trees. A shape predicate selects between the adjacency-only and the renumbering merge strategies, using the renumbering strategy unless the fraction of multi-child internal nodes is below 0.05 (a path-like tree, where renumbering would over-merge). Beyond that, *relaxed* amalgamation [8], applied for  $n \geq 1024$ , merges an adjacent child into its parent even when not fill-justified as long as the merged supernode stays within a width and extra-row budget, trading a little explicit-zero fill for wider, higher-rank Schur updates. Column counts are computed by a depth-first walk of the elimination tree, and drive the factor-size estimate.

## 5 Numeric factorization

### 5.1 Equilibration

Before factoring,  $A$  is symmetrically equilibrated as  $\hat{A} = DAD$  with a real diagonal  $D = \text{diag}(s)$ ,  $s_i = 1/\sqrt{\max_j |A_{ij}|}$ , computed in one pass over the lower triangle; an all-zero row is left unscaled. The real, symmetric, positive diagonal preserves complex symmetry and the pivot signs, tolerates zero diagonals (common in saddle-point systems), and improves the conditioning the bounded pivoting sees. The solve unwinds  $x = D(\hat{A}^{-1}(Db))$ . The LU path applies a two-sided equilibration  $\hat{A} = D_r A D_c$ .

### 5.2 Supernodal left-looking $\text{LDL}^\top$

The default path is a supernodal left-looking  $\text{LDL}^\top$ . Each supernode holds a dense panel of size `nrow`  $\times$  `ncol`, where `nrow` is the number of rows in its factored column structure. The panel is assembled from  $A$ , then updated by every previously factored descendant through the `cmod` operation, which pulls the descendant's columns, applies its block-diagonal  $D$ , and subtracts the resulting rank update; the update runs as a scalar triple loop below a flop threshold and as a SIMD GEMM above it. The fully-summed block is then factored in place by `cddiv`, a blocked Bunch-Kaufman partial factorization with panel width `panel_nb` = 64.

The Bunch-Kaufman kernel [9] factors the dense block as  $P^\top A P = \text{LDL}^\top$  with  $L$  unit lower triangular and  $D$  block-diagonal of  $1 \times 1$  and  $2 \times 2$  pivots. For complex-symmetric  $A$  the control flow is that of the real symmetric `sytf2` with magnitudes  $|z|$  and no conjugation. The pivot threshold is  $\alpha = (1 + \sqrt{17})/8 \approx 0.6404$ ; pivot selection is bounded to the panel's fully-summed rows so the factorization stays within the supernode. A  $2 \times 2$  block is rejected as singular when  $|d_{11}d_{22} - d_{21}^2| < \varepsilon \max(|d_{11}|, |d_{21}|, |d_{22}|)^2$  with  $\varepsilon = 1 \times 10^{-14}$ , which bounds the element growth injected into the trailing update. Each panel's trailing Schur update is deferred and routed through a single SIMD GEMM.

### 5.3 Panel freeing

The transient memory of the left-looking path is controlled by freeing each dense panel once its last consumer is done. Every supernode carries a reference count of the ancestors that will still update from it. When a descendant’s `cmod` decrements that count to zero, the panel is compacted into its final CSC factor fragment and the dense buffer is released; the count is manipulated with acquire/release atomics so this is safe under the tree-parallel factorization. The resident data is therefore the compact factor plus only the live panels, not all panels at once. Sibling subtrees factor concurrently, and the per-node GEMMs parallelize above their flop thresholds.

### 5.4 Multifrontal path

The alternative multifrontal path [10, 11] assembles a dense frontal matrix per supernode from its eliminated columns and the contribution blocks (CBs) of its children, factors the fully-summed part, and passes the Schur complement to the parent as its CB. A work-stealing tree recursion factors children in parallel; each child’s CB is freed the moment it is extend-added into the parent, so only the active contribution frontier is live. Where the child CB stack is large (at least  $4 \times 10^6$  scalars) children are reordered by Liu’s ( $\text{peak} - cb$ ) rule [12] to shrink the transient peak while preserving leaf parallelism; the reordering is a scheduling hint and does not change the numbering, factor, or fill. The multifrontal front is denser and more BLAS-3-rich than a left-looking panel but carries a higher transient peak, which is why the two paths are separate and the choice is per matrix.

### 5.5 Unsymmetric LU

General matrices reuse the symmetric multifrontal machinery on the symmetrized pattern  $A \cup A^\top$ , with UMFPACK-style [13] threshold partial pivoting: within each front’s fully-summed block a diagonal candidate is accepted when  $|a_{jj}| \geq \tau |a_{\text{colmax}}|$  with  $\tau = 0.1$ , else the column maximum is swapped up. The panel width is `NB` = 32, narrower than the  $\text{LDL}^\top$  default because the serial within-panel factorization is the bottleneck while the SIMD GEMM is fast.  $L$  is stored CSC (unit lower),  $U$  CSR; the column permutation is the fill-reducing ordering and a separate row permutation records the pivot interchanges.

### 5.6 Static pivoting

The `ZeroPivotAction` policy governs near-singular pivots. `Fail` (the default) returns a rank-deficiency error on a pivot below the tolerance. `PerturbToEps` lifts any pivot of magnitude below a floor to that floor while preserving its phase,  $d \mapsto d(\text{floor}/|d|)$ , so the factorization satisfies  $\text{LDL}^\top = A + \Delta$  with a bounded  $\Delta$  and never fails; the count of perturbed pivots is tracked per front. This never-fail factor is intended as a preconditioner for iterative refinement; the Krylov path that consumes it is outside the scope of this report.

### 5.7 Kernel scheduling and scoped pools

Four flop-count thresholds route the dense work, passed as a cheap `Copy KernelTuning` value rather than through global state: `scalar_gate` ( $\approx 4096$  flops) below which an update runs as a scalar triple loop instead of a SIMD GEMM; `par_gemm` ( $\approx 1 \times 10^6$ ) above which a `cmod` GEMM runs rayon-parallel; `par_cdiv` ( $\approx 8 \times 10^6$ ) above which the panel-trailing and front GEMMs run parallel, which engages node-level parallelism on the large near-root fronts; and `use_gemm_schur`, which forces the scalar path for benchmarking. The tree recursion runs in a scoped rayon pool of the requested width rather than the global pool, so concurrent solves do not oversubscribe each other. The worker stack is sized to the assembly-tree depth, computed iteratively, `stack = clamp(32 KB · depth, 16 MB, 8 GB)`, which lets deep banded chain-trees factor without a stack overflow.

## 6 A-priori predictors

The analyze phase computes four predictors as exact functions of the symbolic structure and the scalar size  $v = \text{value\_bytes}$ . They are the basis of the memory guarantee, the tuner, the thread policy, and the resource planner.

### 6.1 Per-path peak memory

The factor size is  $\text{factor\_nnz} \cdot (v + 8)$ , where the  $+8$  is the CSC row index and the count is a structural upper bound because numeric cancellation can only lower it. The transient peak is bounded separately for the two paths.

For the left-looking path, a reference-count simulation of the panel-freeing schedule gives the live-panel peak `panel_live_peak_bytes`: it walks the postorder, adds each panel as it is assembled, and frees a panel when its `refcount` reaches zero, recording the running maximum of live panels plus already-compacted factor. This is the floor the solver actually operates at. A conservative whole-run bound assumes the parallel frontier holds all panels at once,

$$\text{scratch} = \frac{1}{4}(\text{panels\_all} + \text{factor}) + 32 \text{ MB}, \quad (2)$$

$$\text{transient} = \text{panels\_all} + \text{factor} + \text{input} + \text{scratch}, \quad (3)$$

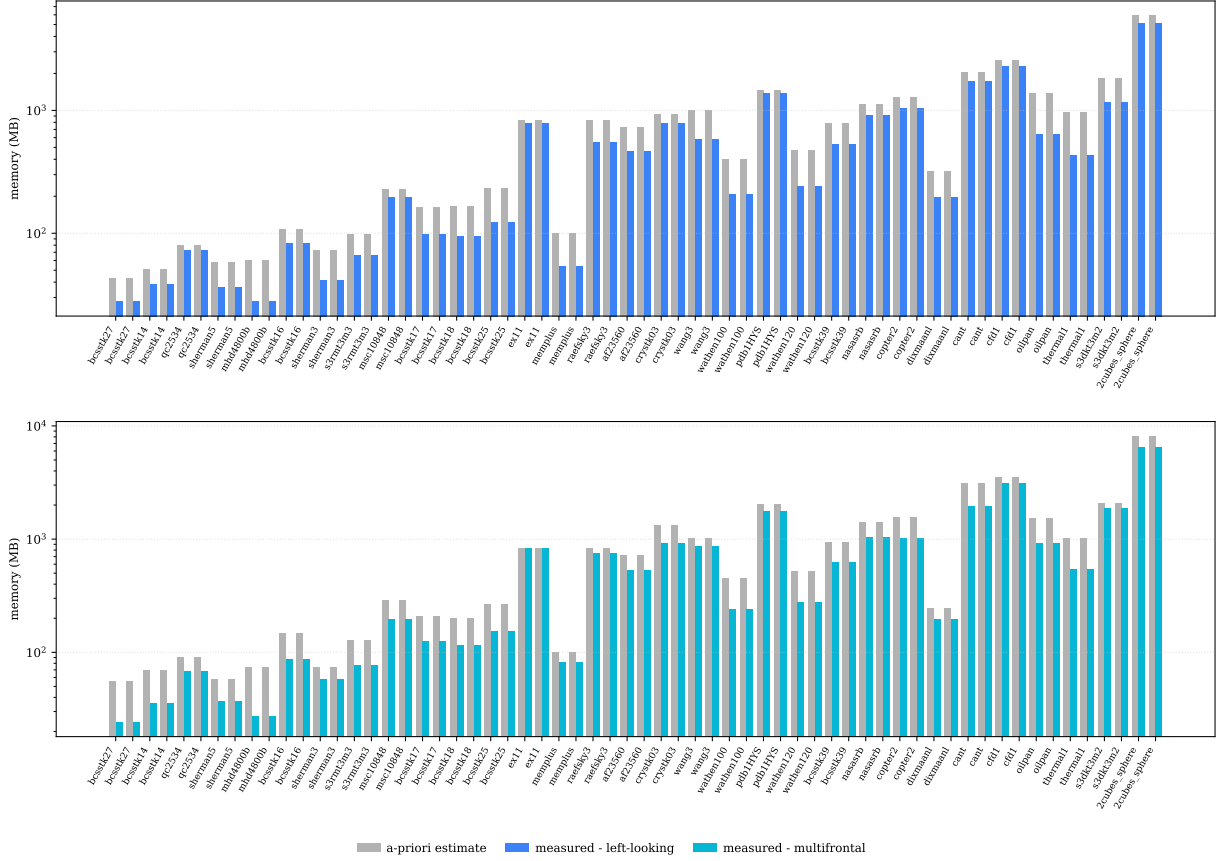
where the scratch term covers the per-thread `cmod/cdiv` and emit buffers. For  $\mathbf{LDL}^\top$  the panel size is  $\text{nrow} \cdot \text{ncol} \cdot v$ , the compacted factor is  $(\frac{nc(nc+1)}{2} + \text{cnrow} \cdot nc)(v + 8)$  per supernode, and the input is  $\text{nnz}(v + 8)$ ; for LU both  $L$  and  $U$  panels are held and the input is  $2 \text{nnz}(v + 8)$ .

For the multifrontal path a level-parallel model takes, at each assembly-tree level, the sum of the front sizes  $\sum \text{nrow}^2 v$  plus the contribution blocks of that level's children  $\sum \text{cnrow}^2 v$ , and multiplies the peak by  $\frac{7}{5}$  to cover the work-stealing overlap of deep leaves of one subtree with mid-level fronts of another. Over the corpus both bounds hold with an estimate/measured ratio of about 1.5 in geomean and no under-prediction (Fig. 2), which is the property the memory guarantee requires; the panel-freeing floor is the tighter quantity used inside the tuner's veto.

### 6.2 Work and runtime

The geometric work  $\text{factor\_flops} = \sum_s \text{nrow}_s^2 \text{ncol}_s$  is a type-independent proxy, so a single `f64` calibration serves all scalar types. The thread-aware runtime estimate is  $\max(\text{critical\_path\_flops}/r, \text{factor\_flops}/(r \cdot \text{speedup}))$  with  $r = \text{gflops} \cdot 10^9$ : the second term is the parallel work, the first is an Amdahl floor from the assembly tree's longest leaf-to-root chain (computed in one postorder pass over the supernodes), so a critical-path-bound matrix does not benefit from workers past the point where the parallel term drops to the floor. The `Calibration` measures the single-thread geometric-flop rate (separately for real and complex, since the per-flop cost differs) and the parallel speedup once, by factoring a  $24^3$  seven-point Laplacian at one thread and at all physical cores, and caches the result keyed by a hardware fingerprint (core counts and a RAM-size bucket) in the manner of FFTW wisdom [14]. The speedup is interpolated linearly from one thread to the calibrated thread count and held flat beyond it, matching the saturating scaling of sparse-direct work; a fallback of 2 Gflop/s and  $\sqrt{\text{cores}}$  speedup is used if the measurement fails.

**Learned residual.** The bare analytical speedup  $s_{\text{pred}}(t) = \text{factor\_flops} / \max(\text{critical\_path\_flops}, \text{factor\_flops}/\text{speedup})$  mispredicts the achieved scaling systematically: the critical-path flops overstate the true serial fraction (sibling subtrees overlap in time), so the model is too pessimistic on serial-heavy trees and too optimistic on wide ones. A ridge regression of  $\log(s_{\text{meas}}/s_{\text{pred}})$  on dimensionless structural features  $[1, \text{amdahl\_frac}, \ln t, \ln \text{tree\_width}]$ , where  $\text{amdahl\_frac} = \text{critical\_path\_flops}/\text{factor\_flops}$ , is fit offline over a generated thread-scaling corpus of 54 matrices and reduces the held-out speedup error by  $\sim 26\%$  under leave-one-matrix-out cross-validation. Almost all of the signal is in `amdahl_frac` (dropping it leaves a 2% gain), which validates the critical-path feature. The residual

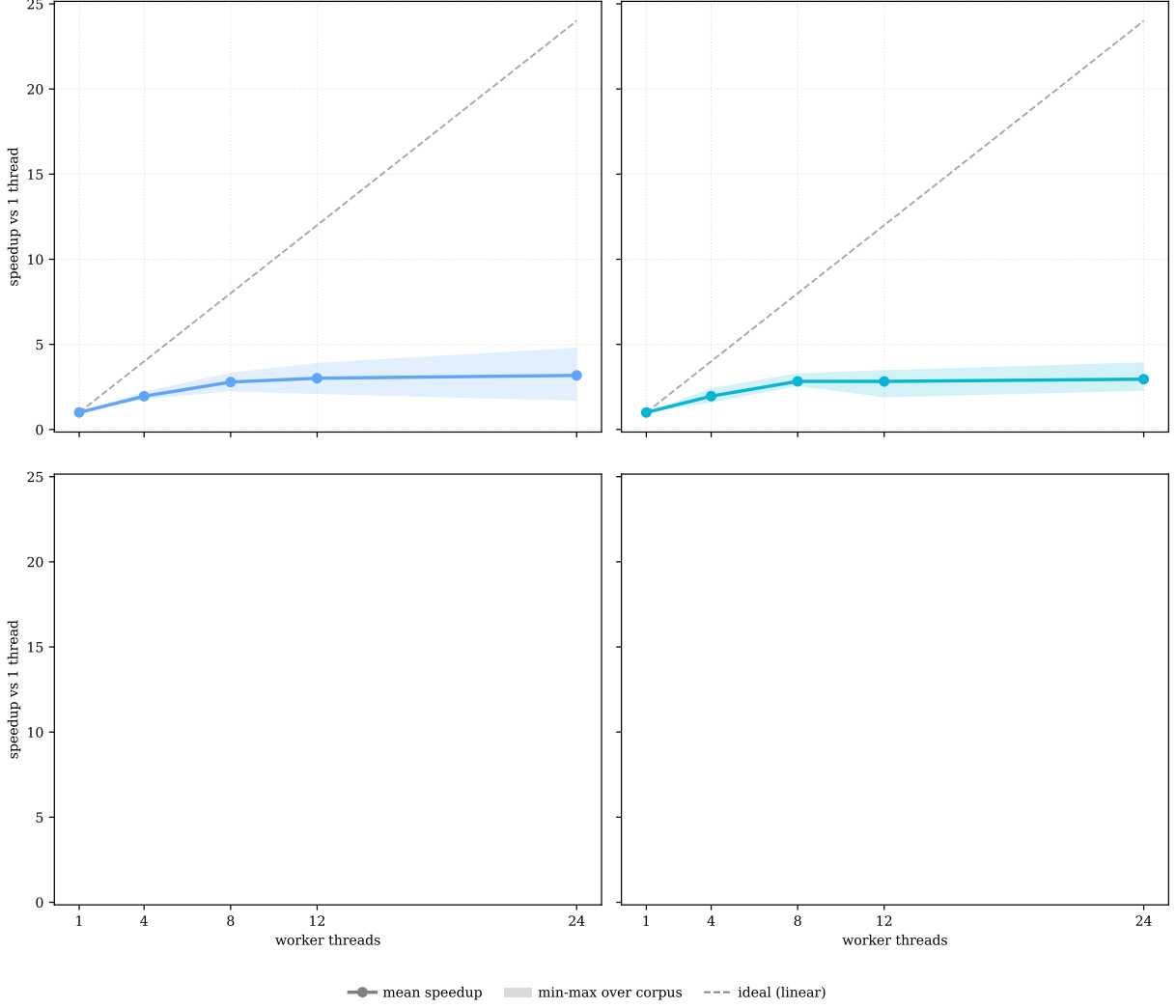


**Figure 2:** A-priori memory estimate (conservative bound and panel-freeing floor) vs. the measured peak, per RSLAB path. The estimate does not under-predict.

is kept additive on the calibrated base rather than a standalone model, so the absolute rate still comes from the hardware calibration and only the shape of the speedup curve is relearned — a machine-transferable correction. It is clamped, and the corrected time is floored at the critical path, so the linear fit can never extrapolate a true chain ( $\text{amdahl\_frac} \rightarrow 1$ ) into a physically impossible speedup.

### 6.3 Thread count

The thread-count predictor maps three structural quantities, the factor flops, the largest front height, and the peak tree width, to a worker count. A thin and narrow problem ( $\text{front\_nrow\_max} < 512$  and  $\text{tree\_width\_max} < 128$ ) is capped at two workers, where more only regress; a small problem ( $\text{factor\_flops} < 3 \times 10^8$ ) is capped at four, where scheduling overhead dominates; otherwise all cores are used. On the corpus this lands within about 10% of the per-matrix optimal in geomean, against about 50% for a fixed budget of two workers. The need for a per-matrix policy is visible in the measured thread scaling (Fig. 3): on the complex corpus RSLAB reaches a geomean  $\sim 3.2\times$  (left-looking) /  $\sim 3.0\times$  (multifrontal) at 12 to 24 workers, but sparse-direct factorization is largely memory-bandwidth bound, so it saturates well below linear, and the min-max band is wide (some matrices regress past a few threads, others reach  $\sim 4.8\times$ ). This spread is why the worker count is picked per matrix — and, in the v2 cost model, a decision variable that may choose fewer threads when scaling saturates.



**Figure 3:** Thread scaling of the RSLAB left-looking and multifrontal kernels, 1–24 workers, geomean with min–max band over the complex corpus.

## 6.4 Structural features

The tuner consumes a structural fingerprint computed in the analyze phase. The pattern features, computed in  $\mathcal{O}(\text{nnz})$ , are  $n$ ,  $\text{nnz}$ , the mean and maximum degree, the degree coefficient of variation (an arrow-pattern signal), the maximum and relative-mean bandwidth, and the diagonal-dominance and diagonal-presence fractions. The symbolic features, computed after analysis, are the supernode count, the fill  $\text{nnz}$  and fill ratio, the mean supernode width, the maximum front height, the tree depth, the maximum and mean tree width, the factor flops, the arithmetic intensity (flops per stored factor entry), and the fraction of flops in the largest and in the top 1% of fronts. These are the quantities that determine which configuration is fastest and smallest, and are the model’s input.

## 7 Tunable stages

Every analyze, factor, and kernel parameter is exposed through one flat `SolverSettings` struct, so a caller and the tuner address the same interface. Table 2 lists the parameters and the stage each affects. The analyze-stage parameters change fill and front shape; the factor-stage parameters change the transient schedule and the pivoting; the kernel-stage parameters change the scalar/SIMD and serial/parallel boundaries without changing the result. The tuner (Sec. 9)

searches the ordering, `nemin`, relaxation width, `panel_nb`, the kernel thresholds, the method, the LU threshold-pivot tolerance `pivot_u`, the equilibration strategy `scaling`, and the memory mode; static pivoting and the incomplete-factor drop tolerance are set directly by the caller. Two further options are caller- or race-selected rather than searched: the RCM band/profile ordering (a candidate in `AutoRace`) and the right-looking method (an opt-in schedule). Every added axis is exact-path: it changes speed or memory but leaves the factorization exact, and each keeps fill predictable so the a-priori memory backstop (Sec. 9) still holds.

Parameter	Stage	Effect
<code>ordering</code>	analyze	Fill-reducing ordering: AMD, AMF, ND (METIS/Scotch/KaHIP), RCM, or a race (Sec. 3)
<code>nemin</code>	analyze	Supernode amalgamation threshold
<code>relax</code>	analyze	Relaxed amalgamation width / extra-row budget
<code>reorder</code>	analyze	Child-reorder strategy (CB-stack peak vs. leaf parallelism)
<code>scaling</code>	analyze	Equilibration strategy (one-pass $\infty$ -norm, iterative Ruiz, MC64, off)
<code>method</code>	factor	Left-looking, multifrontal, or right-looking
<code>threads</code>	factor	Worker count (fixed, or the auto predictor)
<code>pivot_u</code>	factor	LU threshold partial-pivot tolerance $u \in [0, 1]$ ( $u=0$ : static)
<code>on_zero_pivot</code>	factor	Exact fail or static-pivot floor
<code>drop_tol</code>	factor	Incomplete-factor threshold (preconditioner)
<code>memory</code>	factor	Factor emit / peak-RSS strategy
<code>panel_nb</code>	kernel	Bunch–Kaufman / LU panel width
<code>scalar_gate</code>	kernel	Scalar-loop vs. SIMD-GEMM flop gate
<code>par_gemm</code>	kernel	cmod GEMM parallel gate
<code>par_cdiv</code>	kernel	Front / trailing GEMM parallel gate
<code>use_gemm_schur</code>	kernel	SIMD vs. scalar Schur update

**Table 2:** The flat settings interface. Kernel-stage parameters change the schedule of the dense work but not the numeric result.

## 8 Determinism

Three determinism properties hold. First, the numeric factorization is bit-identical regardless of the worker count: the parallel schedule changes the order in which independent subtrees and GEMMs execute, but not the accumulation within a pivot block, so a factor computed at one thread and at all cores agree to the last bit. This extends to the front-level parallelism: the trailing Schur GEMM and the trailing-block subtraction of a large top-of-tree front are split over disjoint output blocks, so the write set is a partition and the result does not depend on the worker count. The multi-RHS triangular solve is likewise bit-identical whether the right-hand sides are solved serially or split into per-thread column chunks (each column is independent), so the parallel wide-RHS solve is a pure speedup. The thread count is therefore a performance parameter only. Second, the analyze-phase choices (the auto ordering, the auto thread count, the amalgamation strategy) are deterministic functions of the matrix pattern. Third, the resource planner (Sec. 10) is a pure function of the memory estimate, the budget, and the cached calibration, so a plan is reproducible and can be computed before the matrix is factored. The pivot sequence is value-dependent (Bunch–Kaufman and threshold partial pivoting both branch on magnitudes), so re-factoring the same numeric matrix reproduces the same pivots; a sequence of matrices that share a pattern but differ in values may pivot differently unless the pivot order is fixed. Setting `pivot_u=0` does exactly that: it keeps the natural pivot order and skips the pivot search across refactorizations (static pivoting for frequency sweeps and time stepping), with iterative refinement recovering accuracy when the values drift. The calibration is a timing measurement and is not bit-reproducible, but once cached it makes subsequent planning on the same hardware reproducible.

## 9 Learned configuration tuner

The tuner selects a `SolverSettings` from the structural feature vector, with **one model per solver path** (symmetric  $\mathbf{LDL}^\top$  and unsymmetric LU): the paths have different relevant axes (`pivot_u` only affects LU) and different speed/memory profiles, so each is fit and selected independently and the LU grid searches the threshold-pivot tolerance the  $\mathbf{LDL}^\top$  grid pins. Each model is a multilayer perceptron predicting (`log factor_ms`, `log peak_mb`) from the standardized features concatenated with an encoding of a candidate configuration: an ordering one-hot, `nemin`, relaxation width, `panel_nb`, `scalar_gate`, `par_cdiv`, the method, the equilibration-strategy one-hot, the memory mode, and (LU) `pivot_u` — a  $37/38 \rightarrow 64 \rightarrow 32 \rightarrow 2$  network with ReLU hidden layers, trained offline with scikit-learn [15] on a complete-distribution corpus (structured-grid generators across real/complex  $\times$  symmetric/unsymmetric  $\times$  conditioning, including curl-curl Maxwell, Stokes/KKT saddle-point, and convection-diffusion over the grid-Péclet range, plus the SuiteSparse matrices) and exported as JSON weights. Each path is fit on its own class: the LU model on 103 unsymmetric matrices (held-out  $R^2$  0.98 in both time and memory), the  $\mathbf{LDL}^\top$  model on 205 symmetric ones. The input encoding is data-driven: the Rust inference builds the same vector the training exported, and a candidate axis is searched only when the shipped model references it, so adding an axis is a retrain, not a code change. Inference is pure Rust and runs at analyze time over the candidate grid; the outputs are destandardized to milliseconds and megabytes. The selected configuration minimizes  $w \log(\hat{t}) + (1 - w) \log(\hat{m})$ , with the default weight  $w = 0.7$ ;  $w = 1$  and  $w = 0$  give the speed and memory extremes.

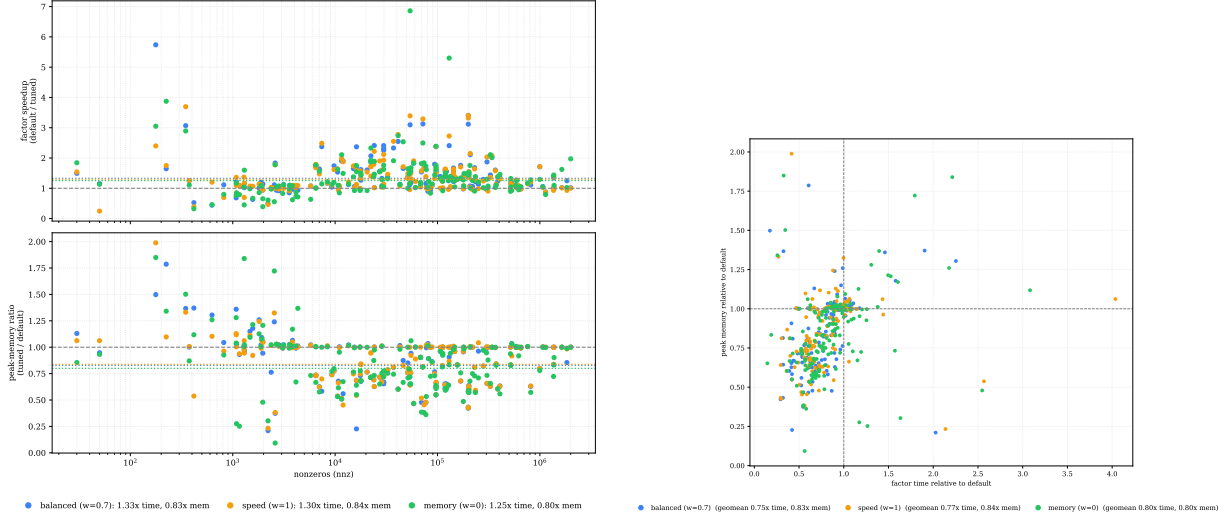
The prediction is constrained by a deterministic guard stack, applied in order:

1. If `factor_flops` exceeds the training range (`flops_ood_cap`  $\approx 1 \times 10^{10}$ ), the model is not trusted; the fallback is not the plain default but a deterministic *ordering race* (`AutoRace`) that picks the minimum-exact-fill ordering (with the default among the candidates, so never worse). The ordering is decidable from the exact a-priori fill, not extrapolated, so this recovers the large-3D / complex-indefinite wins the model cannot: on curl-curl it cut the worst case from  $19\times$  to  $1.5\times$  of PARDISO.
2. A candidate whose predicted `log peak_mb` exceeds the default's by more than  $\ln(1.02)$  is rejected.
3. A multifrontal candidate whose estimated transient exceeds the left-looking floor is rejected, using the deterministic estimate rather than the model.
4. The selected configuration is re-analyzed, and accepted only if its exact estimates satisfy  $\text{fill} \leq 1.02\times$  default,  $\text{flops} \leq 1.05\times$  default, and its memory floor stays under the default's; otherwise the default is used.
5. The survivor must improve the score by at least a deviate threshold (0.08 by default, or 0.20 if it changes the factorization method). Under a hardware-calibrated profile (Sec. 9.1) this threshold is set to  $z \cdot \text{CV}$ , the machine's own single-shot timing noise ( $z = 2$ ), so the tuner never chases a predicted gain smaller than the measurement variance.

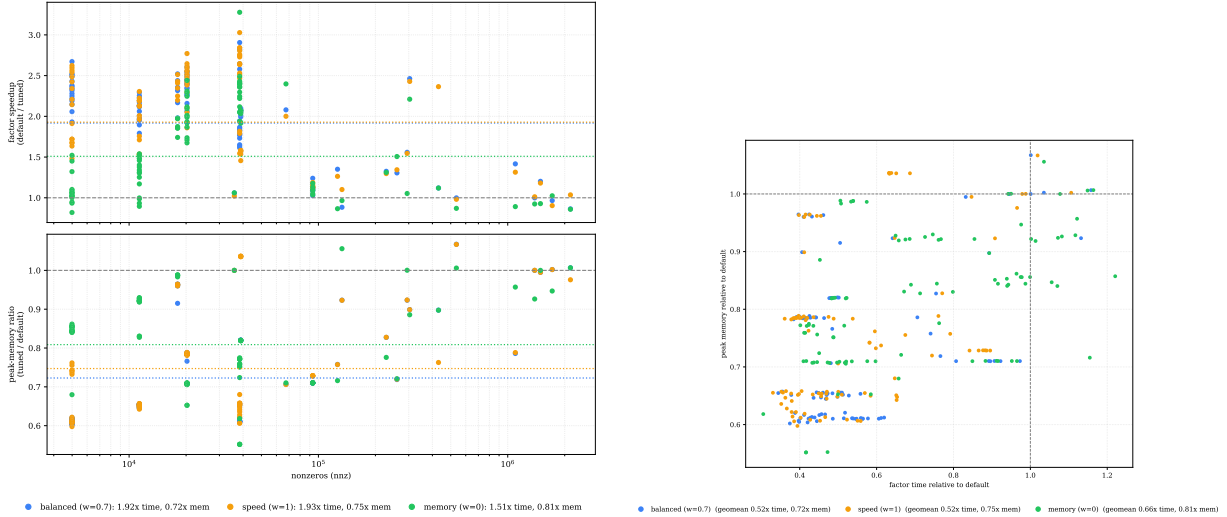
The asymmetry between the memory and time guarantees is a consequence of what is predictable. Peak memory is a deterministic function of structure, so a candidate can be guaranteed to never exceed the default's memory by the estimate-based checks (2)–(4). Factor time depends on the achieved BLAS-3 efficiency, which the model predicts only approximately, so time is optimized in aggregate and backed by the exact-flop proxy in check (4) but is not guaranteed per matrix. The same learned-prediction-with-a-deterministic-fallback structure appears in hardware branch prediction [16], where a misprediction costs only re-execution and not correctness.

Measured against the untuned default, each path is evaluated on *its own* class over 100+ matrices (generated + SuiteSparse), since the two are separate models a caller dispatches to explicitly. The balanced  $\mathbf{LDL}^\top$  tuner ( $w = 0.7$ ) is  $1.33\times$  faster at  $0.83\times$  the memory over 167

matrices; the balanced LU tuner is  $1.92\times$  faster at  $0.72\times$  the memory over 97 matrices, and the guarded picks never regress memory (Figs. 4–5). The LU result is a direct consequence of corpus coverage: with only a handful of generated unsymmetric matrices the LU tuner was statistically indistinguishable from the default ( $1.02\times$ ); broadening the unsymmetric training set to 90 generated convection–diffusion problems (grid-Péclet  $\times$  flow field  $\times$  discretization) lifted the per-path held-out  $R^2$  from 0.75 to 0.98 and the tuner from neutral to  $1.92\times$  (and to  $2.2\times$  on the pure convection–diffusion class). The gains are largest on the big, hard matrices, where a mispicked ordering costs most: on the symmetric head-to-head classes (Sec. 13) the tuner is  $1.82\times$  over the default, its full effect on curl-curl and saddle-point.



**Figure 4:**  $\text{LDL}^T$  tuner vs. untuned default over 167 matrices: speedup and memory ratio by problem size (left), and the three Pareto modes as a time–memory cloud (right).



**Figure 5:** LU tuner vs. untuned default over 97 matrices: speedup and memory ratio by problem size (left), and the three Pareto modes (right). Nearly every matrix lands below-and-left of the default (faster and lighter).

## 9.1 Runtime profile and the meta-tuner

The two per-path models and the guard thresholds ship compiled-in, but they are also exposed as a runtime *tuner profile*: a JSON artifact bundling both models with the calibrated guards. A

profile is applied programmatically (`apply_profile`) or by naming it in the `RSLAB.TUNER.PROFILE` environment variable, which the tuner loads on first use — specializing the solver to a machine or problem class with no recompile. The default profile is a semantic no-op (the embedded models, default guards), so the pathway is exercised on every build.

Profiles are produced by an offline meta-tuner, invoked as `cargo xtask tune`, that chains the whole pipeline deterministically: sweep the corpus, train the two models, measure this machine’s calibration (proxy-flop rate per scalar type, parallel speedup, and the single-shot timing CV that sets the deviate guard), assemble a candidate profile, and *validate* it on a held-out generator corpus (curl-curl and saddle-point at sizes disjoint from training) by factoring each matrix with the shipped default’s recommendation and the candidate’s, and taking the geometric-mean speedup. A *ship-gate* then writes the profile only if it does not regress the default beyond the timing noise floor; a candidate that merely ties still ships (it is the freshly calibrated one), but a regression is rejected and the proven default is kept. This makes retuning safe to run unattended: the worst case a bad sweep or an unlucky machine can produce is the status-quo default, never a slower solver.

## 10 Solver-in-the-loop and resource governance

Two features support running many solves under a fixed resource budget. First, the analyze and factor phases are separate: a symbolic object is computed once and factors any matrix that shares its pattern, so a sequence of same-pattern systems (a frequency sweep, a time-stepping or Newton iteration) pays the ordering and symbolic cost once and only re-runs the numeric factorization. Second, each factorization runs in a scoped thread pool of a caller-set width, so many concurrent solves coexist without oversubscribing the machine; the fixed-budget thread setting is the intended mode for concurrent solving, where the per-solve auto thread count would otherwise contend.

The resource planner `tuning::plan` turns a memory estimate, a budget (a memory ceiling and a thread budget), and the cached calibration into a plan: a predicted peak, a predicted runtime, and a decision on whether the factorization fits the ceiling. It is a pure function of its inputs. Because the memory estimate is an upper bound rather than a sample, a scheduler can pack concurrent solves against a memory ceiling without the safety margin an unbounded estimate would require. Each solve consumes a bounded memory and a thread width; the thread predictor supplies a scaling-aware width per solve, so non-scaling solves can run many-abreast at low thread counts while a BLAS-3-rich solve claims more cores. Because the per-solve peaks and runtimes are predicted, the aggregate memory footprint and the makespan of a batch are computable before the batch runs. The limits on this are the memory-bandwidth contention between concurrent factorizations, which the current calibration does not model (it measures per-solve speedup, not contention), the conservatism of the memory estimate, which reduces packing density, and NUMA placement on multi-socket machines.

## 11 Iterative refinement and preconditioning

The default factorization is exact and the solve is a direct back-substitution. Two modes relax exactness: static pivoting lifts sub-floor pivots so an indefinite or near-singular system never fails, and threshold dropping / block-low-rank make the factor approximate to save memory. In both, the factor is a preconditioner  $M \approx A$  and accuracy is recovered by iteration. RSLAB ships restarted GMRES, a block / multi-RHS GMRES, and the short-recurrence COCG and COCR for the complex-symmetric case, plus classical iterative refinement against the original matrix. The block Krylov variants apply  $M^{-1}$  to all right-hand sides through the parallel `solve_many`, so a block solve inherits its 8 to 19× wide-RHS speedup; and the equilibration and pivot axes (Sec. 7) set the conditioning of that preconditioner factor. The exact-path work of this report is otherwise orthogonal to the Krylov layer — it improves the factor the Krylov layer consumes.

## 12 Test matrix corpus

The tuner is fit and the solver evaluated on a corpus deliberately spanning the full problem distribution — both paths (**LDL<sup>T</sup>** and **LU**), real and complex values, and the conditioning classes (SPD, symmetric indefinite, unsymmetric, ill-scaled, saddle-point) — rather than a single family. The SuiteSparse collection holds few genuinely complex, non-Hermitian matrices, so the backbone is a set of structured-grid generators assembled directly with finite differences (no external FEM library), complemented by the complex SuiteSparse matrices that do exist (the Bai `qc/dwg`, QCD lattice, waveguide, and dielectric-filter EM matrices).

**Curl-curl Maxwell (complex symmetric indefinite).** The time-harmonic operator  $\nabla \times \nabla \times E - (\omega^2 - i\omega\sigma) E$ , discretized on a structured grid through the identity  $\nabla \times \nabla \times = L - \text{grad div}$  (the component-wise vector Laplacian  $L$  minus  $D^T D$  from the discrete divergence  $D$ ), with lumped mass  $M = I$ . The  $-\omega^2 M$  shift makes the operator indefinite; the conductivity term  $i\omega\sigma M$  makes it complex symmetric (not Hermitian); and  $\nabla \times \nabla \times$  is singular on discrete gradients, giving the gradient near-null-space that makes edge-element electromagnetics hard — the primary target class [28, 27].

**Stokes / KKT saddle-point (symmetric indefinite).** The block system  $\begin{bmatrix} A & B^T \\ B & -\beta C \end{bmatrix}$  with  $A$  the velocity vector Laplacian,  $B$  the discrete divergence coupling velocity to pressure, and  $-\beta C$  a Brezzi–Pitkäranta pressure-Laplacian stabilization ( $\beta > 0$ ) regularizing the singular zero (2, 2) block on a collocated grid — the MAC / mixed-FEM form [29, 30].

**Convection–diffusion (unsymmetric).** The operator  $-\varepsilon \nabla^2 u + \mathbf{b} \cdot \nabla u = f$  finite-differenced on a structured grid with Dirichlet boundaries — the canonical unsymmetric class (the LU path’s workload). The first-derivative advection term  $\mathbf{b} \cdot \nabla u$  breaks symmetry; central differencing gives the sharp, oscillation-prone form that stresses pivoting at small  $\varepsilon$  (high grid-Péclet  $|\mathbf{b}|h/\varepsilon$ ), first-order upwinding a diagonally dominant M-matrix. Sweeping  $\varepsilon$  moves the operator from diffusion-dominated (nearly symmetric) to advection-dominated (strongly unsymmetric), and the flow field  $\mathbf{b}$  (a uniform wind or a recirculating vortex) varies the coupling structure. The parameter grid (four Péclet levels  $\times$  two flows  $\times$  two schemes  $\times$  several grids, 2-D and 3-D) yields 90 generated unsymmetric matrices spanning the distribution real CFD/transport solves live on — the LU tuner’s training and held-out corpus.

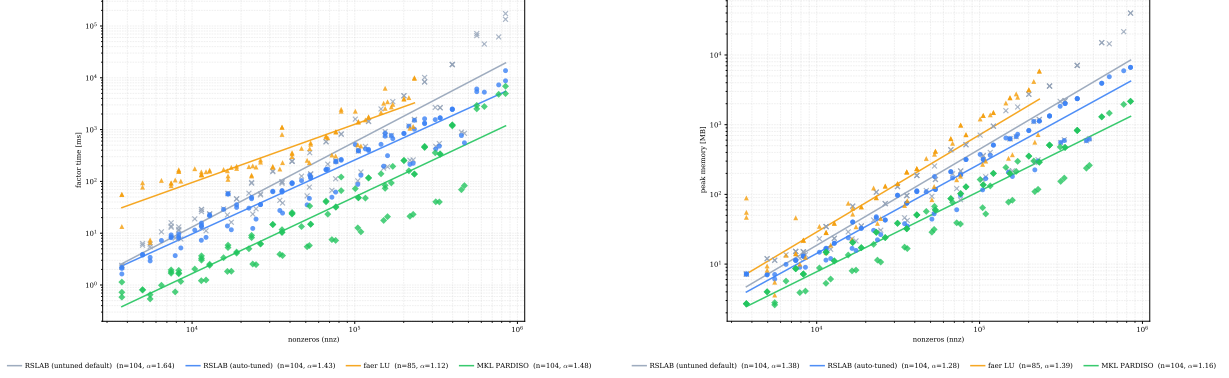
**Remaining classes.** Shifted Helmholtz (complex symmetric), banded / arrow / jumping-coefficient stencils, BEM/MoM near-field kernels (complex unsymmetric), and random SPD / unsymmetric / exactly-conditioned spectral matrices span the rest. Every generator is type-agnostic, so each class is available in `f64` and `Complex<f64>`. This complete-distribution corpus was load-bearing: it exposed the tuner mispicks on complex curl-curl that a real-SPD corpus had hidden, and the too-thin unsymmetric coverage that had left the LU tuner neutral until the convection–diffusion set was added (Sec. 13).

## 13 Evaluation

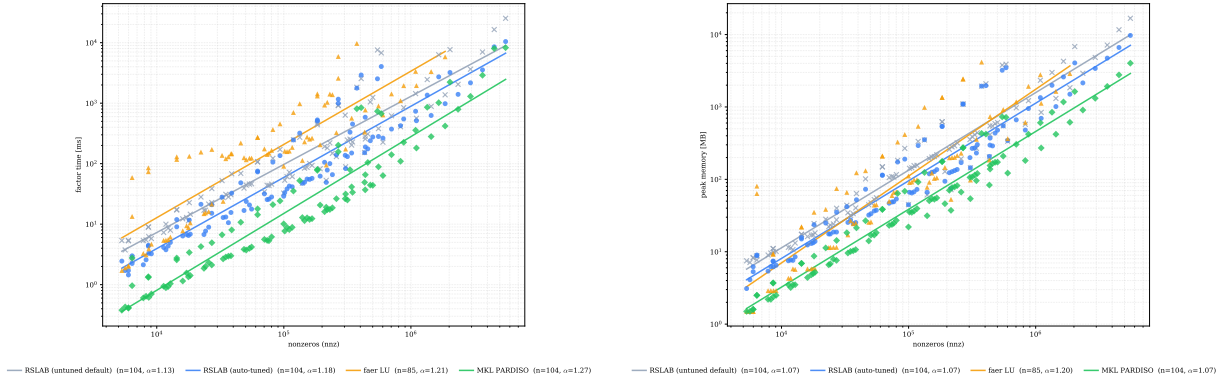
All measured figures come from the `bench_suite` engine over the complete-distribution corpus of Sec. 12 (the structured-grid generators and the complex SuiteSparse [24] matrices, 8k–125k DOFs, all `Complex<f64>`) on a 12-core / 24-thread machine; the per-stage, estimator, thread-scaling, and tuner figures appear with the concepts they support in the preceding sections. This section reports the cross-solver comparison. RSLAB, faer [22], and MKL PARDISO [21] were measured in one run, so the cross-solver numbers are not affected by run-to-run drift.

### 13.1 Scaling

The two paths are separate solvers a caller dispatches to explicitly, so each is compared on *its own* class against its own PARDISO mtype (6 for complex-symmetric, 13 for unsymmetric) and faer: Figs. 6 and 7 plot factor time and peak memory against nnz with a least-squares power-law fit  $\sim \text{nnz}^\alpha$  per solver (solves with residual above 0.1 excluded). The RSLAB curve is the auto-tuned default, which picks left-looking or multifrontal per matrix under the guards. Table 3 gives the head-to-head geomean ratios.



**Figure 6:**  $\text{LDL}^\top$  path (symmetric, PARDISO mtype 6): factor time (left) and peak memory (right) vs. problem size — RSLAB vs. faer vs. MKL PARDISO, with power-law fits.



**Figure 7:** LU path (unsymmetric, PARDISO mtype 13): factor time (left) and peak memory (right) vs. problem size. Each plot carries the untuned-default RSLAB curve (gray) beside the auto-tuned one (blue); the tuner closes part of the default→PARDISO distance. On time RSLAB scales slightly flatter than PARDISO ( $\alpha \approx 1.17$  vs. 1.26).

RSLAB sits between the two on both paths: faster and lighter than the pure-Rust faer, moderately behind the hand-optimized MKL PARDISO. The unsymmetric path is the closer of the two ( $3.9\times$  vs.  $5.8\times$  behind PARDISO). Each head-to-head plot carries the untuned default curve beside the tuned one, so the learned tuner’s win (last two table rows,  $1.82\times/1.65\times$  faster at  $0.68\times/0.71\times$  the memory) is visible as a gap that widens with problem size — the tuner helps most on the big matrices where a mispicked ordering costs most, and never at a memory cost: the backstop compares the *exact* symbolic fill (not a dense-panel estimate, which overshoots non-uniformly across orderings and once let a  $2\times$ -fill pick slip through), so a pick can never grow the factor beyond the default’s. The corpus choice was load-bearing throughout: benchmarking on real-SPD matrices flattered RSLAB ( $7\times$ ) and hid two tuner mispicks on complex-indefinite curl-curl, where the adaptive Auto ordering produces  $\sim 3\times$  the fill of nested dissection. Adding MetisND to the tuner grid and racing orderings out-of-distribution (the AutoRace fallback, deciding on the exact a-priori fill rather than the extrapolating model) cut the worst curl-curl case from

RSLAB vs	$\mathbf{LDL}^\top$ (sym)	LU (unsym)
MKL PARDISO — factor time	$5.8\times$ slower	<b><math>3.9\times</math></b> slower
MKL PARDISO — peak memory	$2.3\times$ more	$2.4\times$ more
faer $\text{LU}^\dagger$ — factor time	<b><math>6.9\times</math></b> faster	<b><math>3.7\times</math></b> faster
faer $\text{LU}^\dagger$ — peak memory	<b><math>2.3\times</math></b> less	$1.1\times$ less
untuned default — factor time	<b><math>1.82\times</math></b> faster	<b><math>1.65\times</math></b> faster
untuned default — peak memory	<b><math>0.68\times</math></b> (less)	<b><math>0.71\times</math></b> (less)

**Table 3:** Head-to-head geomean ratios (104 matrices per path, 5k–1M nonzeros), each path on its own class ( $\mathbf{LDL}^\top$ : curl-curl, Helmholtz, saddle-point; LU: convection–diffusion, BEM/MoM), over the matrices both solvers factor to below 0.1 residual. The last row is the auto-tuned solver’s speedup over the untuned RSLAB default.  $^\dagger$ faer has no symmetric path (it factors symmetric matrices as LU), so its  $\mathbf{LDL}^\top$  gap is structurally largest; it also OOMs on the largest matrices (factoring only the smaller subset), so its head-to-head is a conservative floor.

$19\times$  to  $1.5\times$  of PARDISO.

## 13.2 Accuracy

Figure 8 shows the relative residual  $\|Ax - b\|/\|b\|$  per solver. Where RSLAB factors, 24 of 31 matrices are below  $1 \times 10^{-8}$ , matching PARDISO and ahead of faer. The exact  $\mathbf{LDL}^\top$  declines a handful of indefinite saddle-point matrices rather than returning a degraded solution.

## 13.3 Exact-path axes: correctness and speedups

The tunable axes were validated over 180 SuiteSparse matrices spanning structural, CFD, thermal, and saddle-point classes. Every determinism and equivalence property held on all 180: the right-looking factor is bit-identical to the multifrontal one; the **Eager** and **LowMemory** emit modes produce bit-identical factors; the parallel front subtraction is bit-identical to serial; and the 32-bit compressed factor solves bit-identically to the full-width one. Where the axes carry a measurable gain, the parallel multi-RHS solve is 8 to  $19\times$  faster than solving the 256 right-hand sides column by column, and 32-bit index compression halves the stored index footprint (for example 459 MB  $\rightarrow$  230 MB on `pwtk`) at no accuracy cost. The block/multi-RHS Krylov solvers (block GMRES, COCG, COCR) apply the factored preconditioner through exactly this parallel `solve_many`, so they inherit the wide-RHS speedup directly, and the equilibration/pivot axes set the quality of that preconditioner factor; the exact-path work is otherwise orthogonal to the Krylov layer. Retraining the tuner with the new axes is neutral on the well-scaled structural corpus (geomean  $1.03\times$  versus the pre-axis model, within the  $\pm 20\%$  single-shot variance): the threshold-pivot and equilibration axes bind on unsymmetric and ill-scaled classes, and the memory backstop keeps the guarded picks from regressing where they do not help.

## 14 Related work

Automatic and learned tuning of sparse solvers has been studied along several axes. GPTune [17] tunes HPC libraries including SuperLU\_DIST by Bayesian optimization with multitask learning across problem sizes, using online trial runs per problem; the RSLAB model is a feed-forward regressor evaluated once at analyze time. OSKI [18] autotunes at the sparse-kernel level (register and cache blocking for sparse matrix-vector products), not at the solver-configuration level. Lighthouse [19] and the work of Bhowmick, Raghavan et al. [20] apply machine learning to select a solver or preconditioner from a portfolio, a classification problem, rather than to predict and tune a direct factorization’s configuration. Production direct solvers such as MKL PARDISO [21] and MUMPS ship hand-written heuristic automatic modes rather than learned models. The structure used here, a learned performance regressor constrained by a deterministic guarantee

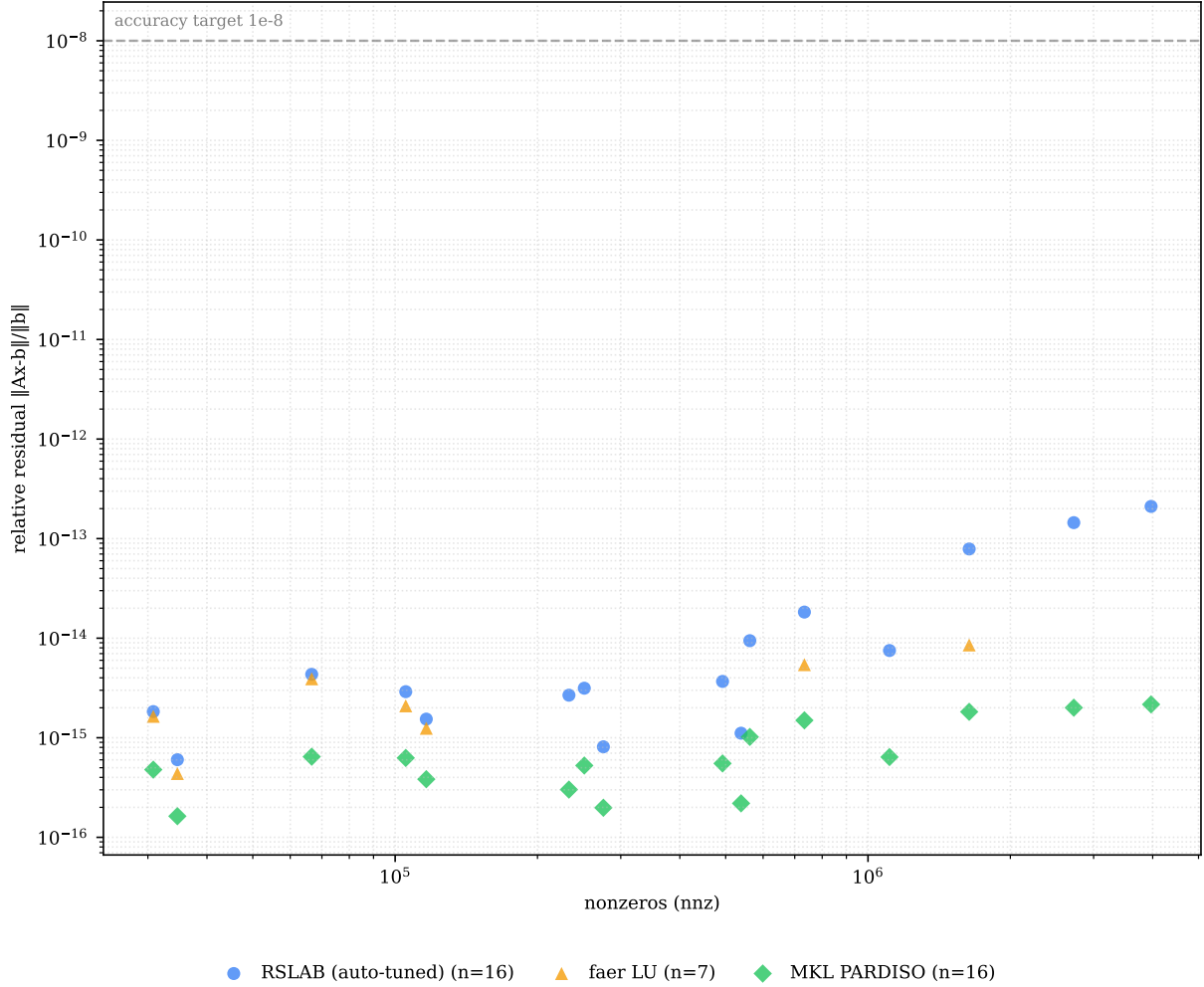


Figure 8: Relative residual across the corpus, per solver.

and an out-of-range fallback, is closest in form to hardware branch prediction [16], where a learned predictor is paired with a mechanism that makes a misprediction cost re-execution and not correctness.

## 15 Limitations and future work

The tuner’s gains are aggregate: the geomean speedup is  $1.62\times$ , and part of it is within the per-matrix single-shot timing variance of about  $\pm 20\%$ , so the geomean rather than any single matrix is the signal. The model is trained on one machine and one corpus and does not transfer across hardware without retraining, in the same way a branch predictor is tuned per microarchitecture; a meta-tuner that runs the sweep, the training, and the guard-threshold calibration on a user’s own matrix set to emit a problem-class profile is planned. Time is optimized but not guaranteed; only memory is. The exact-path tuning stages that keep the factorization exact are now implemented (Sec. 7): an RCM/band ordering [25], an exposed threshold-pivot tolerance, a tunable equilibration strategy [26], static pivot-order reuse for fixed-pattern sequences via `pivot.u=0`, a parallel blocked multi-RHS solve, a right-looking schedule, and front-level parallelism of the trailing Schur update and subtraction. Each provides an a-priori memory estimate so the deterministic memory guarantee continues to hold, and each either enters the tuner grid or is caller-selected. What remains is a fuller two-dimensional (block-cyclic) factorization of the very largest top-of-tree fronts to push thread scaling past the current  $\sim 3\times$ , and the cross-hardware meta-tuner above.

## 16 Conclusion

RSLAB is a pure-Rust, scalar-generic sparse direct solver that reaches within a single-digit factor of MKL PARDISO on the SuiteSparse corpus and is faster and smaller than the open-source alternatives. Its analyze phase computes a per-path memory bound, a work and runtime estimate, and a thread recommendation, all as exact functions of structure; a learned tuner selects the solver configuration from that structure under a deterministic guard stack that keeps its memory below the untuned default. The numeric result is independent of the worker count and the resource planner is a pure function of its inputs, which makes batched, solver-in-the-loop execution reproducible and its cost computable in advance.

## References

- [1] P. R. Amestoy, T. A. Davis, I. S. Duff. An approximate minimum degree ordering algorithm. *SIAM J. Matrix Anal. Appl.* 17(4):886–905, 1996; and Algorithm 837: AMD. *ACM Trans. Math. Softw.* 30(3):381–388, 2004.
- [2] P. R. Amestoy. Recent progress in parallel multifrontal solvers and the approximate minimum fill (HAMF) ordering. Technical report / CERFACS, 1999.
- [3] E. Rothberg, S. C. Eisenstat. Node selection strategies for bottom-up sparse matrix ordering. *SIAM J. Matrix Anal. Appl.* 19(3):682–695, 1998.
- [4] A. George. Nested dissection of a regular finite element mesh. *SIAM J. Numer. Anal.* 10(2):345–363, 1973.
- [5] G. Karypis, V. Kumar. A fast and high quality multilevel scheme for partitioning irregular graphs. *SIAM J. Sci. Comput.* 20(1):359–392, 1998.
- [6] C. M. Fiduccia, R. M. Mattheyses. A linear-time heuristic for improving network partitions. *19th Design Automation Conference*, 175–181, 1982.
- [7] J. D. Hogg, E. Ovtchinnikov, J. A. Scott. A sparse symmetric indefinite direct solver for GPU architectures (SSIDS). *ACM Trans. Math. Softw.* 42(1):1–25, 2016.
- [8] C. Ashcraft, R. Grimes. The influence of relaxed supernode partitions on the multifrontal method. *ACM Trans. Math. Softw.* 15(4):291–309, 1989.
- [9] J. R. Bunch, L. Kaufman. Some stable methods for calculating inertia and solving symmetric linear systems. *Math. Comp.* 31(137):163–179, 1977.
- [10] I. S. Duff, J. K. Reid. The multifrontal solution of indefinite sparse symmetric linear equations. *ACM Trans. Math. Softw.* 9(3):302–325, 1983.
- [11] J. W. H. Liu. The multifrontal method for sparse matrix solution: theory and practice. *SIAM Review* 34(1):82–109, 1992.
- [12] J. W. H. Liu. On the storage requirement in the out-of-core multifrontal method for sparse factorization. *ACM Trans. Math. Softw.* 12(3):249–264, 1986.
- [13] T. A. Davis. Algorithm 832: UMFPACK, an unsymmetric-pattern multifrontal method. *ACM Trans. Math. Softw.* 30(2):196–199, 2004.
- [14] M. Frigo, S. G. Johnson. The design and implementation of FFTW3. *Proc. IEEE* 93(2):216–231, 2005.
- [15] F. Pedregosa et al. Scikit-learn: machine learning in Python. *J. Mach. Learn. Res.* 12:2825–2830, 2011.
- [16] D. A. Jiménez, C. Lin. Dynamic branch prediction with perceptrons. *7th Int. Symp. High-Performance Computer Architecture (HPCA)*, 197–206, 2001.

- [17] Y. Liu, W. M. Sid-Lakhdar, O. Marques, X. Zhu, C. Meng, J. W. Demmel, X. S. Li. GPTune: multitask learning for autotuning exascale applications. *PPoPP*, 234–246, 2021.
- [18] R. Vuduc, J. W. Demmel, K. A. Yelick. OSKI: a library of automatically tuned sparse matrix kernels. *J. Phys.: Conf. Ser.* 16:521–530, 2005.
- [19] E. R. Jessup, P. Motter, B. Norris, K. Sood. Performance-based numerical solver selection in the Lighthouse framework. *SIAM J. Sci. Comput.* 38(5):S750–S771, 2016.
- [20] S. Bhowmick, V. Eijkhout, Y. Freund, E. Fuentes, D. Keyes. Application of machine learning in selecting sparse linear solvers. *Int. J. High Perf. Comput. Appl.*, 2006.
- [21] O. Schenk, K. Gärtner. Solving unsymmetric sparse systems of linear equations with PARDISO. *Future Gener. Comput. Syst.* 20(3):475–487, 2004.
- [22] S. Quinones. faer-rs: a linear algebra library for the Rust programming language. <https://github.com/sarah-quinones/faer-rs>, 2024.
- [23] X. S. Li. An overview of SuperLU: algorithms, implementation, and user interface. *ACM Trans. Math. Softw.* 31(3):302–325, 2005.
- [24] T. A. Davis, Y. Hu. The University of Florida sparse matrix collection. *ACM Trans. Math. Softw.* 38(1):1–25, 2011.
- [25] E. Cuthill, J. McKee. Reducing the bandwidth of sparse symmetric matrices. *Proc. 24th ACM National Conference*, 157–172, 1969.
- [26] D. Ruiz. A scaling algorithm to equilibrate both rows and columns norms in matrices. Technical Report RAL-TR-2001-034, Rutherford Appleton Laboratory, 2001.
- [27] J.-C. Nédélec. Mixed finite elements in  $\mathbb{R}^3$ . *Numer. Math.* 35(3):315–341, 1980.
- [28] J.-M. Jin. *The Finite Element Method in Electromagnetics*, 3rd ed. Wiley–IEEE Press, 2014.
- [29] H. C. Elman, A. Ramage, D. J. Silvester. Algorithm 866: IFISS, a Matlab toolbox for modelling incompressible flow. *ACM Trans. Math. Softw.* 33(2):14, 2007.
- [30] M. Benzi, G. H. Golub, J. Liesen. Numerical solution of saddle point problems. *Acta Numerica* 14:1–137, 2005.