

Bryan Van Scoy

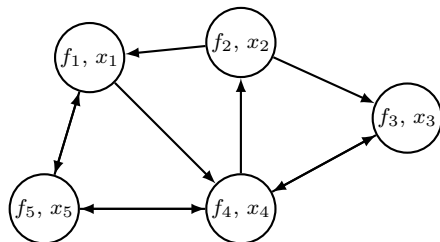
✉ bvanscoy@miamioh.edu 🌐 [vanscoy.github.io](https://github.com/vanscoy)

Miami University

Distributed optimization

$$\underset{x_1, \dots, x_n}{\text{minimize}} \quad \sum_{i=1}^n f_i(x_i)$$

$$\text{subject to} \quad x_1 = x_2 = \dots = x_n$$



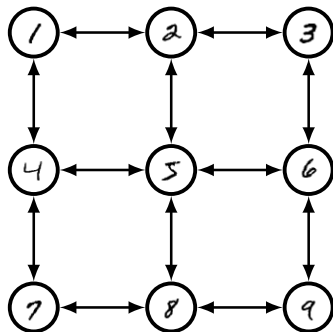
$$W = \begin{bmatrix} \frac{1}{3} & \frac{1}{3} & 0 & 0 & \frac{1}{3} \\ 0 & \frac{1}{2} & 0 & \frac{1}{2} & 0 \\ 0 & \frac{1}{6} & \frac{2}{3} & \frac{1}{6} & 0 \\ \frac{1}{3} & 0 & \frac{1}{3} & 0 & \frac{1}{3} \\ \frac{1}{3} & 0 & 0 & \frac{1}{3} & \frac{1}{3} \end{bmatrix}$$

Want each agent to compute the global optimizer x_* by communicating with local neighbors and performing local computations.

Application: Distributed machine learning

Each agent i has

- data (x_i, y_i)
- features $\Phi(x_i)$
- prediction model $\theta_i^T \Phi(\cdot)$
- error $\|y_i - \theta_i^T \Phi(x_i)\|$



$$\begin{aligned} & \underset{\theta_1, \dots, \theta_n}{\text{minimize}} && \sum_{i=1}^n \|y_i - \theta_i^T \Phi(x_i)\| \\ & \text{subject to} && \theta_1 = \theta_2 = \dots = \theta_n \end{aligned}$$

Distributed gradient descent (DGD)

$$x_i^+ = \sum_{j=1}^n w_{ij} x_j - \alpha \nabla f_i(x_i)$$

- linear convergence to *suboptimal* solution with constant stepsize
- *sublinear* convergence to optimal solution with decaying stepsize

The optimal solution is *not* a fixed point.

Distributed inexact gradient tracking (DIGing)

$$x_i^+ = \sum_{j=1}^n w_{ij} x_j - \alpha y_i$$

$$y_i^+ = \sum_{j=1}^n w_{ij} y_j + \nabla f_i(x_i^+) - \nabla f_i(x_i)$$

- linear convergence over time-varying networks
- the bound grows with the number of agents
- the bound does not apply to large stepsizes

Have conservative bounds that are specific to DIGing.

Other distributed algorithms

$$x^+ = Wx - \alpha \nabla f(x) \quad \text{DGD'09}$$

$$x^{++} = (I + W)x^+ - \frac{I+W}{2}x - \alpha (\nabla f(x^+) - \nabla f(x)) \quad \text{EXTRA'15}$$

$$x^{++} = W(2x^+ - Wx - \alpha W(\nabla f(x^+) - \nabla f(x))) \quad \text{AugDGM'15}$$

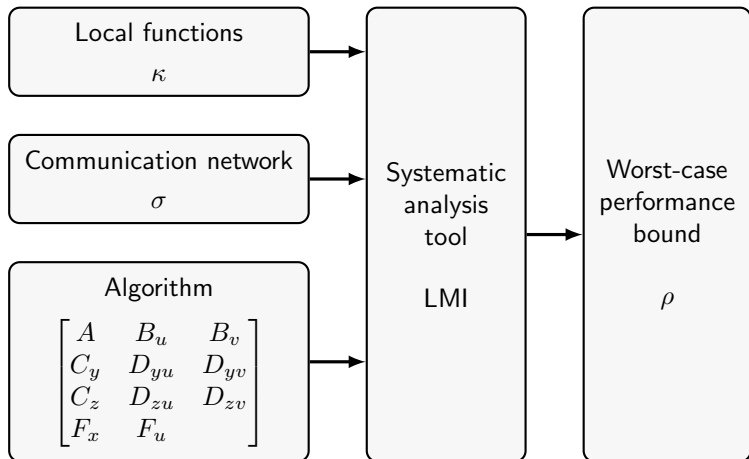
$$x^{++} = W(2x^+ - Wx) - \alpha (\nabla f(x^+) - \nabla f(x)) \quad \text{DIGing'17}$$

$$x^{++} = (I + W)x^+ - \frac{I+W}{2}x - \alpha (\nabla f(\frac{I+W}{2}x^+) - \nabla f(\frac{I+W}{2}x)) \quad \text{ExDiff'17}$$

$$x^{++} = (I + W)x^+ - \frac{I+W}{2}(x + \alpha (\nabla f(x^+) - \nabla f(x))) \quad \text{NIDS'19}$$

and more...

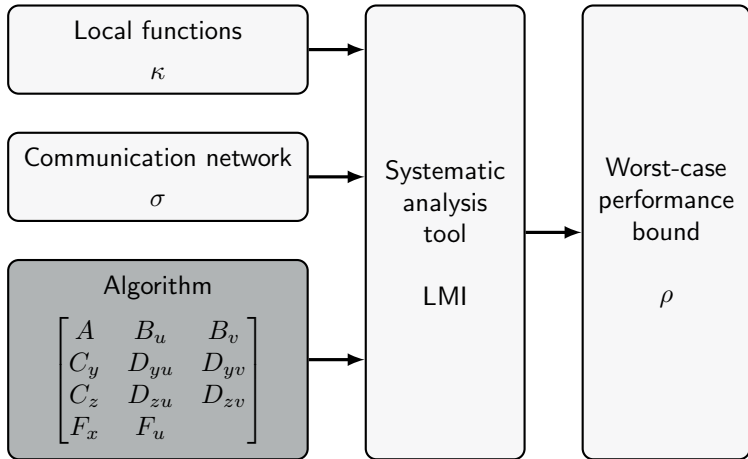
A systematic approach

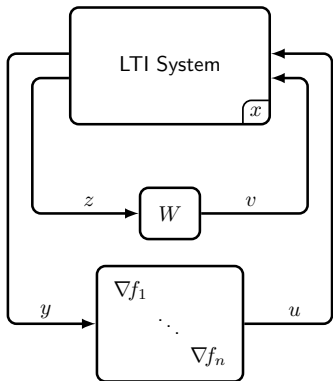


Main idea

$$V(x) = \underbrace{\|\text{avg}(x) - x_\star\|_P^2}_{\text{optimality}} + \underbrace{\|x - \text{avg}(x)\|_Q^2}_{\text{consensus}}$$

- search for a quadratic Lyapunov function using LMIs
- use pointwise quadratic constraints from the assumptions on the local objective functions and the communication network
- similar to Lur'e problem from robust control





$$\begin{bmatrix} x_i^+ \\ y_i \\ z_i \end{bmatrix} = \begin{bmatrix} A & B_u & B_v \\ C_y & D_{yu} & D_{yv} \\ C_z & D_{zu} & D_{zv} \end{bmatrix} \begin{bmatrix} x_i \\ u_i \\ v_i \end{bmatrix}$$

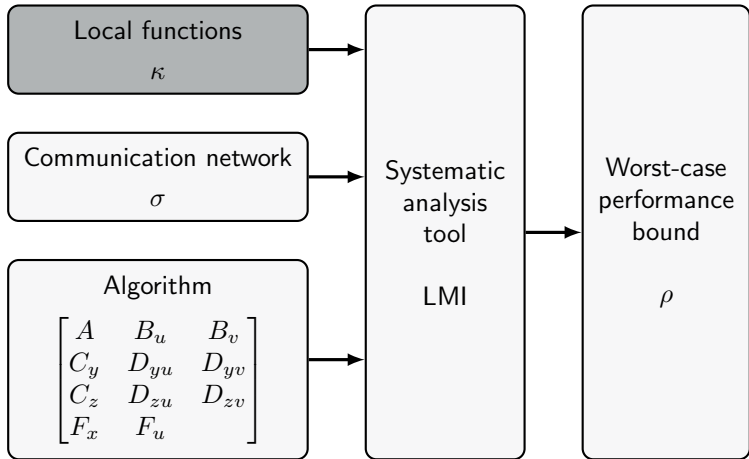
$$v_i = \sum_{j=1}^n w_{ij} z_j$$

$$u_i = \nabla f_i(y_i)$$

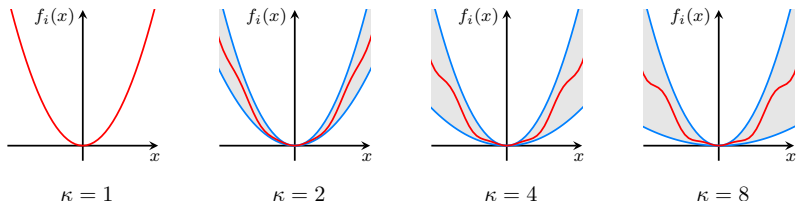
$$0 = \sum_{j=1}^n (F_x x_j + F_u u_j)$$

Many known algorithms fit in this form.

	2 state variables	3 state variables
1 communicated variable	SVL $\begin{bmatrix} 1-\gamma & \beta & -\alpha & -\gamma \\ -1 & 1 & 0 & -1 \\ 1-\delta & 0 & 0 & \delta \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}$	EXTRA $\begin{bmatrix} 1 & -\frac{1}{2} & \alpha & -\alpha & 1 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 0 \\ 1 & -\frac{1}{2} & 0 & 0 & 0 \\ 1 & -1 & \alpha & 0 & 0 \end{bmatrix}$
	Exact Diffusion (ExDIFF) $\begin{bmatrix} 1 & -1 & -\alpha & 1 \\ \frac{1}{2} & 0 & -\alpha & \frac{1}{2} \\ 1 & 0 & -\frac{1}{2} & 0 \\ 1 & 0 & 0 & 0 \\ 1 & -1 & 0 & 0 \end{bmatrix}$	NIDS $\begin{bmatrix} 1 & -\frac{1}{2} & \frac{\alpha}{2} & -\frac{\alpha}{2} & 1 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 0 \\ 1 & -\frac{1}{2} & \frac{\alpha}{2} & -\frac{\alpha}{2} & 0 \\ 1 & -1 & \alpha & 0 & 0 \end{bmatrix}$
2 communicated variables	Unified DIGing (uDIG) $\begin{bmatrix} 0 & -\alpha & -\alpha & 1 & 0 \\ \frac{L+m}{2} & 0 & -1 & 0 & 1 \\ 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 \\ -\frac{L+m}{2} & 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \end{bmatrix}$	DIGing $\begin{bmatrix} 0 & -\alpha & 0 & 0 & 1 & 0 \\ 0 & 0 & -1 & 1 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & -\alpha & 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & -1 & 0 & 0 & 0 \end{bmatrix}$
	Unified EXTRA (uEXTRA) $\begin{bmatrix} 0 & -\alpha & -\alpha & 1 & 0 \\ 0 & 0 & -1 & L & 1 \\ 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & -L & 0 \\ 0 & 1 & 0 & 0 & 0 \end{bmatrix}$	AugDGM $\begin{bmatrix} 0 & 0 & 0 & 0 & 1 & -\alpha \\ 0 & 0 & -1 & 1 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & -\alpha \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & -1 & 0 & 0 & 0 \end{bmatrix}$

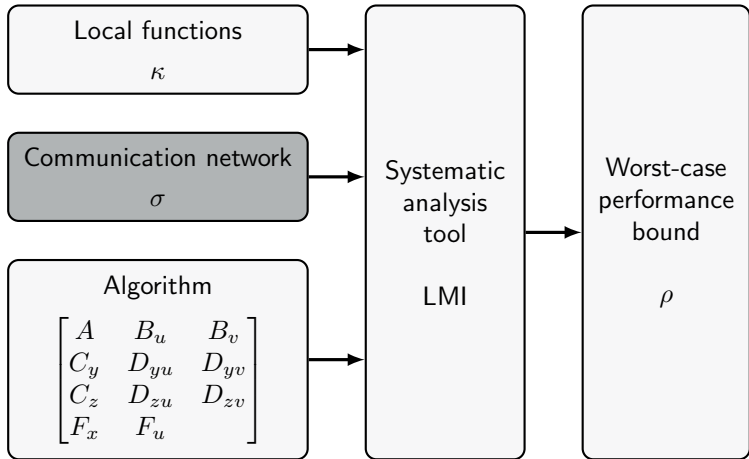


Each local objective function f_i is L -smooth and m -strongly convex with respect to the global optimizer x_* .

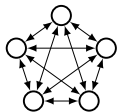


The condition ratio $\kappa = L/m$ characterizes the variation in curvature.

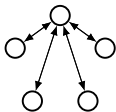
$$\begin{bmatrix} \nabla f_i(y) - \nabla f_i(y_*) \\ y - y_* \end{bmatrix}^\top \begin{bmatrix} -2mL & L + m \\ L + m & -2 \end{bmatrix} \begin{bmatrix} \nabla f_i(y) - \nabla f_i(y_*) \\ y - y_* \end{bmatrix} \geq 0$$



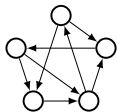
At each iteration, the gossip matrix W satisfies the sparsity pattern of the graph, is doubly stochastic, and has spectral gap σ .



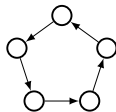
$$\sigma = 0.0$$



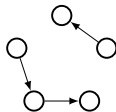
$$\sigma = 0.667$$



$$\sigma = 0.707$$



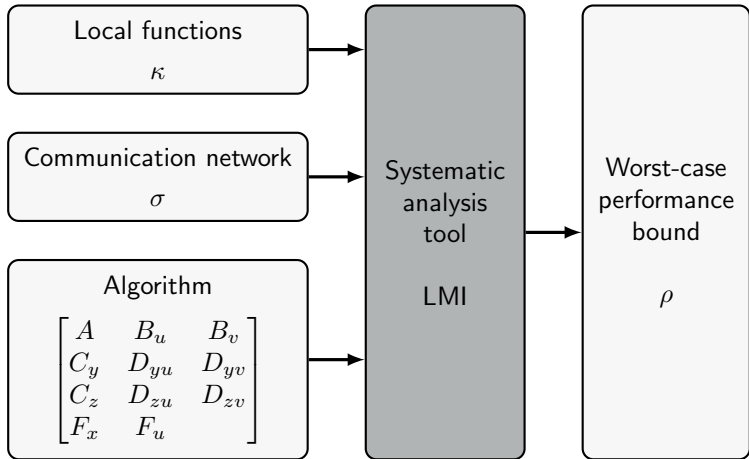
$$\sigma = 0.809$$



$$\sigma = 1.0$$

The spectral gap $\sigma = \left\| \frac{1}{n} \mathbf{1}\mathbf{1}^T - W \right\|_2$ characterizes network connectivity.

$$\begin{bmatrix} z - \text{avg}(z) \\ Wz - \text{avg}(z) \end{bmatrix}^T \begin{bmatrix} \sigma^2 & 0 \\ 0 & -1 \end{bmatrix} \begin{bmatrix} z - \text{avg}(z) \\ Wz - \text{avg}(z) \end{bmatrix} \geq 0$$



Optimality:

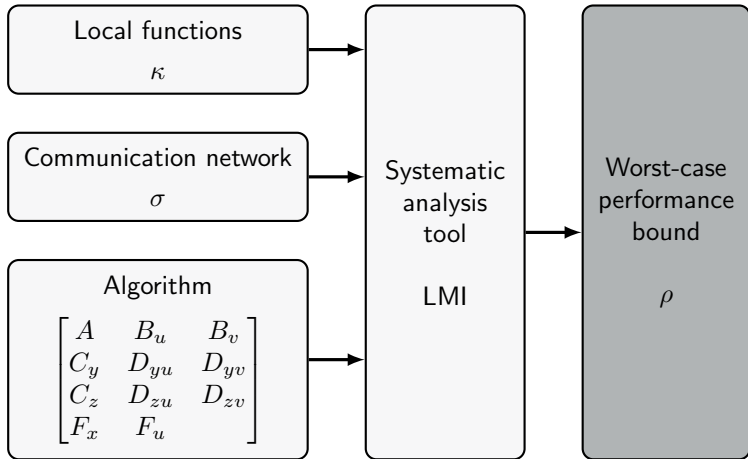
$$\Psi^T \begin{bmatrix} A & B_u \\ I & 0 \\ \hline C_y & D_{yu} \\ 0 & I \end{bmatrix}^T \begin{bmatrix} P & 0 & | & 0 \\ 0 & -\rho^2 P & | & 0 \\ \hline 0 & 0 & | & M_0 \end{bmatrix} \begin{bmatrix} A & B_u \\ I & 0 \\ \hline C_y & D_{yu} \\ 0 & I \end{bmatrix} \Psi \preceq 0$$

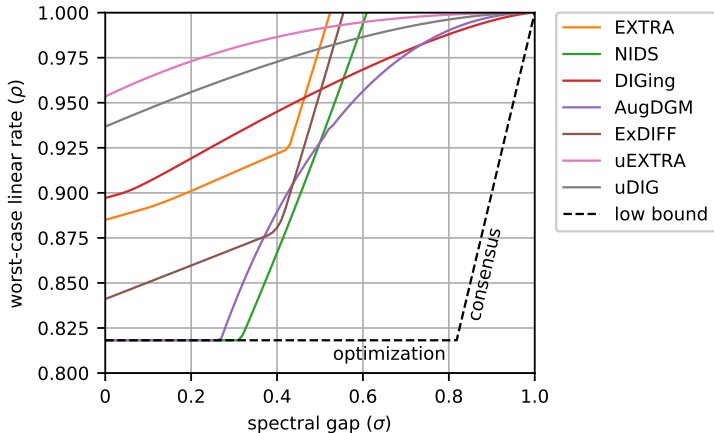
Consensus:

$$\begin{bmatrix} A & B_u & B_v \\ I & 0 & 0 \\ \hline C_y & D_{yu} & D_{yv} \\ 0 & I & 0 \\ \hline C_z & D_{zu} & D_{zv} \\ 0 & 0 & I \end{bmatrix}^T \begin{bmatrix} Q & 0 & | & 0 & | & 0 \\ 0 & -\rho^2 Q & | & 0 & | & 0 \\ \hline 0 & 0 & | & M_0 & | & 0 \\ 0 & 0 & | & 0 & | & M_1 \otimes R \end{bmatrix} \begin{bmatrix} A & B_u & B_v \\ I & 0 & 0 \\ \hline C_y & D_{yu} & D_{yv} \\ 0 & I & 0 \\ \hline C_z & D_{zu} & D_{zv} \\ 0 & 0 & I \end{bmatrix} \preceq 0$$

$$M_0 = \begin{bmatrix} -2\kappa & \kappa + 1 \\ \kappa + 1 & -2 \end{bmatrix} \quad M_1 = \begin{bmatrix} \sigma^2 - 1 & 1 \\ 1 & -1 \end{bmatrix} \quad \Psi = \text{null} [F_x \quad F_u]$$

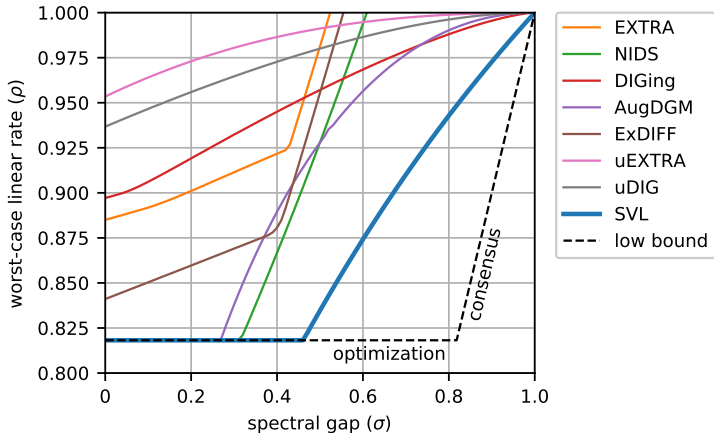
Feasibility implies each agent converges to x_* linearly with rate ρ .





- $\kappa = 10$ and the gradient stepsize α is optimized
- lower bound corresponds to only optimization and only consensus

Can analyze *all* algorithms using the same technique.



- SVL has the best worst-case convergence
- equivalent to inexact ADMM (**Boyd et al., 2011**)

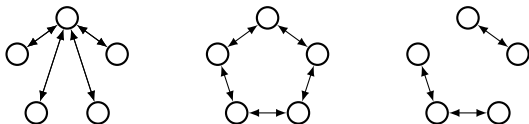
The analysis LMIs can also be used for *design*.

Extensions

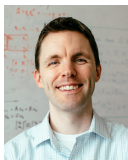
- optimized trade-off between communication and computation

iteration	1	2	3	4	5	6	7	8	9
communication	✓	✓	✓	✓	✓	✓	✓	✓	✓
gradient evaluation			✓			✓			✓

- jointly-connected communication networks



the network is connected only when averaged over time



Laurent Lessard



Akhil Sundararajan

TCNS'20 A. Sundararajan, B. Van Scoy, and L. Lessard, "Analysis and design of first-order distributed optimization algorithms over time-varying graphs"

CDC'20 B. Van Scoy and L. Lessard, "Systematic analysis of distributed optimization algorithms over jointly-connected networks"

ACC'19 A. Sundararajan, B. Van Scoy, and L. Lessard, "A canonical form for first-order distributed optimization algorithms"

NecSys'19 B. Van Scoy and L. Lessard, "A distributed optimization algorithm over time-varying graphs with efficient gradient evaluations"

<https://vanscoy.github.io>