

# Acceleration of Newton-like methods for nonlinear systems by preflattening techniques

Congrès National d'Analyse Numérique

June 2, 2026

Ngoc Do Quyen DANG

Supervisors: Quang Huy TRAN, Clément CANCÈS, Ibtihel BEN GHARBIA



# Outline

## 1 Context and objectives

- Motivations
- Towards a novel approach

## 2 Notion of flatness and preflattening

## 3 Preflattening: direct approaches

## 4 Preflattening: elimination approaches

## 5 Conclusion and perspectives

# Need for a faster nonlinear solver

- Challenge: solving **large nonlinear** algebraic systems derived from the numerical discretization of physical models
- Newton's method
  - Pros:
    - **high** convergence rate (superlinear or even quadratic) with the *sufficiently close* initials
  - Cons:
    - local convergence
    - slow convergence or failure: strong and unbalanced nonlinearities
- Goal: improving the effectiveness of Newton's method
- Existing techniques:
  - better starting points
  - globalization, line search, trust region methods
  - continuation and homotopy methods
  - nonlinear preconditioning

# Nonlinear preconditioning

## ■ Nonlinear system

$$F(\mathbf{u}) = 0$$

<i>Left preconditioning</i>	<i>Right preconditioning</i>
$G(F(\mathbf{u})) = G(0)$	$F(U(\boldsymbol{\tau})) = 0$
- Same unknown variables - <b>Modified</b> equations	- <b>Change</b> unknown variables - Same equations

## ■ Pros:

- ↑ robustness: Bassetto '21, Jonval '24
- ↑ speed: Marelli '25

## ■ Cons:

- Linear systems: lower the *condition number* = faster iterative solvers: CG, GMRES, BiCGSTAB,...
- Nonlinear systems: **no** equivalent quantity that governs convergence rate and that is decreased by preconditioning

# New technique: Preflattening

- Affine function: Newton's method converges in a **single** iteration
- Paradigm: the **closer**  $F$  is to an affine function, the **faster** Newton's method converges
- Basic idea:
  - identifying new quantity:  $\rho_F(\mathbf{u})$  called *flatness number*
  - designing an equivalent system at the current point s.t.

$$\rho_{F \circ U}(U^{-1}(\mathbf{u})) < \rho_F(\mathbf{u})$$

⇒ **preflattening** technique

- Advantage: act exclusively at the nonlinear level of Newton's method
  - disentangling the two sources of difficulty
  - leaving the linear level to existing linear preconditioners

# Outline

- 1 Context and objectives
- 2 Notion of flatness and preflattening**
  - Notion of flatness
  - Flatness and rate of convergence
  - General right-preflattening scheme
- 3 Preflattening: direct approaches
- 4 Preflattening: elimination approaches
- 5 Conclusion and perspectives

# Flatness matrix - Flatness number

- Newton's *iteration function*

$$\Phi_F(\mathbf{u}) = \mathbf{u} - \nabla F(\mathbf{u})^{-1} F(\mathbf{u})$$

- The **flatness matrix** of  $F$  at  $\mathbf{u} \in \Omega$  is the Jacobian matrix of the iteration function

$$\mathcal{J}_F(\mathbf{u}) := \nabla \Phi_F(\mathbf{u}) = \nabla F(\mathbf{u})^{-1} \nabla^2 F(\mathbf{u}) \nabla F(\mathbf{u})^{-1} F(\mathbf{u})$$

- Ostrowski's theorem: *spectral radius* of  $\mathcal{J}_F$  at the solution  $\mathbf{u}^*$  plays a crucial role for convergence
- The **flatness number** of  $F$  at  $\mathbf{u} \in \Omega$  is spectral radius of the flatness matrix

$$\rho_F(\mathbf{u}) := \rho(\mathcal{J}_F(\mathbf{u}))$$

- $\rho_F(\mathbf{u})$ : dimensionless and affine-invariant

# Connection with Newton's rate of convergence

## ■ Kantorovich '48

- semi-local convergence of Newton's method
- **quadratic** convergence rate governed by the product of 2 constants

$$L_K = \sup_u \|\nabla F(\mathbf{u}^0)^{-1} \nabla^2 F(\mathbf{u})\|, \quad \eta = \|\nabla F(\mathbf{u}^0)^{-1} F(\mathbf{u}^0)\|$$

- $\rho_F(\mathbf{u}^0)$ : lower bound of this product

## ■ Hernández and Salanova '94

- similar statements by a single constant  $M = \sup_u \|\mathcal{J}_F(\mathbf{u})\|$
- $\rho_F(\mathbf{u}^0)$ : lower bound of  $M$

## ■ Halley's method

- high-order Newton's scheme with **cubic** convergence rate

$$\mathbf{u}^{k+1} = \mathbf{u}^k - \left[ I - \frac{1}{2} \mathcal{J}_F(\mathbf{u}^k) \right]^{-1} \nabla F(\mathbf{u}^k)^{-1} F(\mathbf{u}^k)$$

- the smaller the correction term  $\frac{1}{2} \mathcal{J}_F(\mathbf{u}^k)$  is, the closer Halley is to Newton

# A new approach: right-preflattening

- $\{U_\alpha\}_{\alpha \in \mathcal{A}}$ : a family of parametrization  $\mathbf{u} = U_\alpha(\boldsymbol{\tau})$
- $\mathcal{A}$  can be **discrete** or **continuous**
- Given the current iterate  $\mathbf{u}^k$ , we successively compute

$$\boldsymbol{\alpha}^k = \arg \min_{\alpha \in \mathcal{A}} \{ \rho_{F \circ U_\alpha}(U_\alpha^{-1}(\mathbf{u}^k)) \}, \quad [\text{optimal preflattening}]$$

$$\boldsymbol{\tau}^k = U_{\alpha^k}^{-1}(\mathbf{u}^k), \quad [\text{current variable}]$$

$$\begin{aligned} \boldsymbol{\tau}^{k+1} &= \boldsymbol{\tau}^k - [\nabla_{\boldsymbol{\tau}}(F \circ U_{\alpha^k})(\boldsymbol{\tau}^k)]^{-1}(F \circ U_{\alpha^k})(\boldsymbol{\tau}^k) \\ &= \boldsymbol{\tau}^k - [\nabla_{\boldsymbol{\tau}} U_{\alpha^k}(\boldsymbol{\tau}^k)]^{-1} [\nabla_{\mathbf{u}} F(\mathbf{u}^k)]^{-1} F(\mathbf{u}^k), \quad [\text{updated variable}] \end{aligned}$$

$$\mathbf{u}^{k+1} = U_{\alpha^k}(\boldsymbol{\tau}^{k+1}), \quad [\text{primary unknown}]$$

- Halley's method = right-preflattener in this framework

# Outline

- 1 Context and objectives
- 2 Notion of flatness and preflattening
- 3 Preflattening: direct approaches**
  - Models
  - Methods
  - Numerical results
- 4 Preflattening: elimination approaches
- 5 Conclusion and perspectives

# Chemical equilibrium problem

- Speciation problem: find the number of moles  $n_j$  of species  $C_j$  s.t.

$$\mathbf{A}\mathbf{n} = \mathbf{f}, \quad [\text{element conservation}]$$

$$S^T \boldsymbol{\mu}(\mathbf{n}) = 0, \quad [\text{equality of chemical potentials}]$$

- Expanded reformulation [Jonval, 2024]

$$f_1 := a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n - b_1 = 0,$$

$$\vdots \quad \quad \quad \vdots \quad \quad \quad \vdots \quad \quad \quad \vdots \quad \quad \quad \vdots \quad \quad \quad \vdots$$

$$f_e := a_{e1}x_1 + a_{e2}x_2 + \dots + a_{en}x_n - b_e = 0,$$

$$f_{e+1} := s_{11} \ln x_1 + s_{12} \ln x_2 + \dots + s_{1n} \ln x_n - t_1 = 0,$$

$$\vdots \quad \quad \quad \vdots \quad \quad \quad \vdots \quad \quad \quad \vdots \quad \quad \quad \vdots \quad \quad \quad \vdots$$

$$f_n := s_{\tau 1} \ln x_1 + s_{\tau 2} \ln x_2 + \dots + s_{\tau n} \ln x_n - t_\tau = 0$$

- Difficulty: as  $x_i \rightarrow 0^+$ , the derivatives of  $\ln x_i$  **blow up**

# Equivalent formulations

- The **log-trick** model  $G(\mathbf{y}) = 0$  is defined by replacing

$$\ln x_i \leftarrow y_i \quad \text{and} \quad x_i \leftarrow \exp y_i$$

- Log-trick algorithm: Newton's method on  $G$
- Drawback: **explosion** of derivatives of  $\exp y_i$  as  $y_i \rightarrow +\infty$
- The **parameterized** model  $H(\boldsymbol{\tau}) = 0$  is established by replacing

$$\ln x_i \leftarrow Y_i(\tau_i) \quad \text{and} \quad x_i \leftarrow X_i(\tau_i)$$

- Parametrization  $(X_i(\tau_i), Y_i(\tau_i))$  depends on each method
  - Enhance robustness
  - Speed up convergence

# Brenner-Cancès (BC) method

- Objective: guarantee robustness in all solution regimes
- BC parametrization

$$(X_i(\tau_i), Y_i(\tau_i)) = \begin{cases} (\exp \tau_i, \tau_i) & \text{if } \tau_i < 0 \\ (\tau_i + 1, \ln(\tau_i + 1)) & \text{if } \tau_i \geq 0 \end{cases}$$

- Procedure: Start from  $(x^0, y^0)$ ,
  - compute parameter  $\tau^0$
  - run Newton on  $H$  to update  $\tau$  until the stopping criteria hold
- **Robust** by construction

# Continuous Dynamic (CD) method

- Motivation: Using a **larger** family of parametrizations to better minimize the flatness number
- CD parametrization

$$X_{\alpha_i}(\tau_i) = (1 + (1 - \alpha_i)\tau_i)^{\frac{1}{1-\alpha_i}}$$
$$Y_{\alpha_i}(\tau_i) = \frac{1}{1 - \alpha_i} \ln(1 + (1 - \alpha_i)\tau_i)$$

- Minimize an **upper bound**  $\hat{\rho}_{H_\alpha} := \left\| P^{-1} \mathcal{J}_{H_\alpha}(\boldsymbol{\tau}) P \right\|_{\text{Frob}}$  to find optimal values for  $\alpha_i$
- Numerical observations:
  - BC is robust but sometimes **slower** than CD
  - CD is less robust but sometimes **faster** than BC

# Hybrid method

- Aim: taking advantage of BC's **robustness** and CD's **local rapidity**
- Idea: use CD only when the upper bound of the flatness number  $\widehat{\rho}_{H_\alpha}$  is *small enough*
- Procedure: At each iteration  $k$ ,
  - Given  $(x_i^k, y_i^k)$  with  $y_i^k = \ln(x_i^k)$
  - Compute  $\widehat{\rho}_{H_\alpha}^k$
  - Run BC or CD:
    - If  $\widehat{\rho}_{H_\alpha}^k > \alpha$ , run BC,  
where  $\alpha$  is prescribed threshold, e.g.,  $\alpha = 0.5$
    - Otherwise, run CD

# Numerical data

## ■ Test cases

- H<sub>2</sub>O: 3 species, 2 elements, 1 reaction
- NaCl: 8 species, 4 elements, 4 reactions
- Seawater: 37 species, 10 elements, 27 reactions

## ■ Stopping conditions

$$\|H(\tau)\| < \varepsilon = 10^{-10}, \quad \text{max\_iter} = 100$$

## ■ High-precision stopping conditions

$$\|H(\tau)\| < \varepsilon = 10^{-150}, \quad \text{max\_iter} = 100$$

## ■ Methods

- **Log-trick**: a state-of-the-art technique in industry
- **BC**: a technique preventing the blow-up of derivatives
- **CD**: a technique that is integrated with preflattening
- **Hybrid**: a combination of BC and CD methods

# Results of NaCl case

## ■ Number of iterations

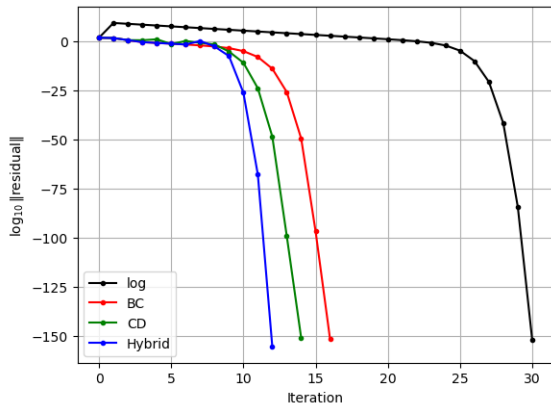
Initial guess		log	BC	CD	Hybrid
$n_1^0$	$n_{j \neq 1}^0$				
55	55	26	12	13	10
55	1	14	12	11	10
55	0.1	11	11	11	10
55	$10^{-2}$	12	12	12	11
55	$10^{-4}$	x	30	x	27
55	$\varepsilon_{32}$	x	16	x	16

## ■ Remarks

- CD is less robust but sometimes faster than BC
- Hybrid is the **fastest** algorithm

# Evolution of residuals in NaCl case

## High-precision figure



## CD and Hybrid methods

- Speed up Newton's convergence
- Not produce results as impressive as those of the scalar counterpart

# Outline

- 1 Context and objectives
- 2 Notion of flatness and preflattening
- 3 Preflattening: direct approaches
- 4 Preflattening: elimination approaches**
  - Elimination as a preconditioning technique
  - Numerical results
- 5 Conclusion and perspectives

# Elimination strategy

- Problem: **unbalanced nonlinearities**  $\Rightarrow$  Newton's slowness
- Idea: eliminate “bad” unknowns from corresponding “bad” equations
- Splitting of good/bad variables:  $\tau = (\tau_g, \tau_b)$  and equations

$$H(\tau) = 0 \Leftrightarrow \begin{cases} H_g(\tau_g, \tau_b) = 0 \\ H_b(\tau_g, \tau_b) = 0 \end{cases}$$

- Elimination scheme
  - From bad equations  $H_b = 0$ , extract bad unknowns as functions of good ones, i.e.,  $\tilde{\tau}_b = \varphi_b(\tau_g)$
  - Apply Newton to the reduced system  $H_g(\tau_g, \varphi_b(\tau_g))$
- Difficulty: **complicated calculation** due to a Schur complement

# Equivalent two-step procedure

- *2-step formulation*: at each iteration  $k$

- **Projection**: solve for  $\tilde{\tau}_b^k = \varphi_b(\tau_g)$  the *inner* system

$$H_b(\tau_g^k, \tilde{\tau}_b^k) = 0$$

and set the projected state as

$$\tilde{\tau}^k = \varphi(\tau^k) := (\tau_g^k, \tilde{\tau}_b^k) = (\tau_g^k, \varphi_b(\tau_g^k))$$

- **Linearization**: perform one global *outer* Newton iteration starting from  $\tilde{\tau}^k$ , i.e.,

$$\tau^{k+1} = \tilde{\tau}^k - \nabla H(\tilde{\tau}^k)^{-1} H(\tilde{\tau}^k)$$

- **Advantage**: dynamic partitioning into good/bad subset [Brenner et al. 2026]
- **Challenge**:
  - PDEs: obvious to associate a “bad” unknown with a “bad” function
  - Chemical systems: no notion of mesh

# Connection between equations and unknowns

- **Linear** normalized residues

$$(\xi_1, \xi_2, \dots, \xi_n)^T = \Xi(\bar{\mathbf{y}}; \boldsymbol{\tau}) := \nabla G(\bar{\mathbf{y}})^{-1} H(\boldsymbol{\tau}),$$

where Jacobian matrix of log-trick formulation  $\nabla G$  is computed at the current state  $\bar{\mathbf{y}}$

- Benefits of considering  $\Delta(\boldsymbol{\tau}) = 0$ 
  - Equivalent to  $H(\boldsymbol{\tau}) = 0$
  - Same orders of magnitude as those of unknowns
  - Easy to link with “bad” unknowns: choose the same indices for bad equations and bad unknowns

# How to choose bad equations?

- Techniques inspired from statistics
  - Clustering: k-means, IQR, MAD...
  - Heuristic criteria: median, max distance...
- Novel approach based on flatness
  - Upper bound

$$\widehat{\rho}_{H_\alpha} = \|P^{-1} \mathcal{J}_{H_\alpha}(\tau) P\|_{\text{Frob}} = \left( \sum_{j=1}^n \Gamma_j \right)^{1/2}$$

where

$$\Gamma_j = \|[\mathcal{J}_{H_\alpha}]_{j\text{-th column}}\|_2^2$$

- Idea: cancel those columns  $j$  that have **largest** norms  $\Gamma_j$

$$\mathcal{B} = \{j \mid \Gamma_j \text{ most significantly contributes to } \widehat{\rho}_{H_\alpha}\}$$

# Numerical tests

## ■ Test cases

- H<sub>2</sub>O: 3 species, 2 elements, 1 reaction
- NaCl: 8 species, 4 elements, 4 reactions
- Seawater: 37 species, 10 elements, 27 reactions

## ■ Stopping conditions

$$\|H(\tau)\| < \varepsilon = 10^{-10}, \quad \text{max\_iter} = 100$$

## ■ Methods

- **BC**: a technique proposed in Jonval's thesis to enhance the robustness
- **Hybrid**: a combination of BC and CD methods
- **CDe**: an integration of elimination into the CD method

# The necessity of preflattening

- H<sub>2</sub>O test case
- Algorithms integrated with elimination
  - CDe: elimination + CD method (with preflattening)
  - BCe: elimination + BC method (without preflattening)
- Number of iterations

Initial guess		BC	Hybrid	BCe	CDe
$n_2^0$	$n_3^0$				
$10^{-2}$	$2 \cdot 10^{-2}$	7	4	4	2
$10^{-2}$	$5 \cdot 10^{-2}$	7	4	4	2
$10^{-4}$	$2 \cdot 10^{-4}$	7	4	4	2
$10^{-4}$	$5 \cdot 10^{-4}$	7	4	4	2
$\varepsilon_{32}$	$2\varepsilon_{32}$	7	4	4	2
$\varepsilon_{32}$	$5\varepsilon_{32}$	7	4	4	2

- CDe is a better than BCe
  - ⇒ **optimizing** the flatness number is essential

# Seawater test case

## Complexity

$$m^2 \times (\text{inner iterations}) + n^2 \times (\text{outer iterations})$$

## Number of iterations

Initial guess		BC	Hybrid	CDe
$n_1^0$	$n_{j \neq 1}^0$			
55	55	22	22	<b>7</b>
55	1	19	16	<b>6</b>
55	0.1	17	14	<b>5</b>
55	$10^{-2}$	16	14	<b>5</b>
55	$10^{-4}$	27	27	<b>5</b>
55	$\varepsilon_{32}$	29	27	<b>5</b>

Number of outer iterations

Initial guess		BC	Hybrid	CDe
$n_1^0$	$n_{j \neq 1}^0$			
55	55	31768	31768	<b>21921</b>
55	1	27436	23104	<b>19322</b>
55	0.1	24548	20216	<b>16206</b>
55	$10^{-2}$	23104	20216	<b>13787</b>
55	$10^{-4}$	38988	38988	<b>13743</b>
55	$\varepsilon_{32}$	41876	38988	<b>24693</b>

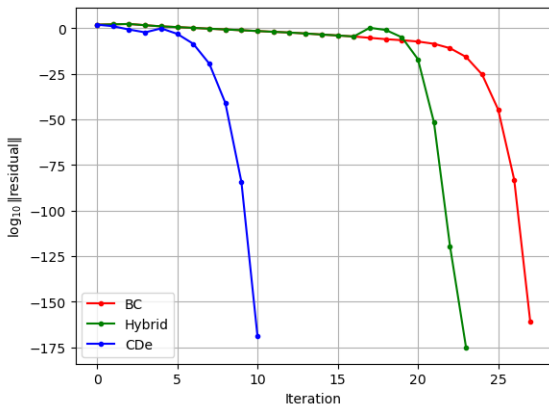
Complexity

## Remarks

- CDe outperforms the others
- CDe requires the fewest iterations, always less than 10

# Evolution of residuals in Seawater case

## High-precision figure



- BC and Hybrid experience a quite long plateau before decreasing
- CDe is the fastest algorithm

# Outline

- 1 Context and objectives
- 2 Notion of flatness and preflattening
- 3 Preflattening: direct approaches
- 4 Preflattening: elimination approaches
- 5 Conclusion and perspectives**

# Key results and future works

## ■ Conclusions

- Flatness number is a promising measure of the nonlinearity of a function
- Preflattening is a technique used to reduce the flatness number
- Elimination helps accelerate the convergence rate of Newton's method, while minimizing the flatness number is also indispensable

## ■ Future works

- Discretization of PDEs in geosciences: Richards equation
- Theoretical framework for the new approach

**Thank you for listening**