



Background & Motivation

Summary

- We propose a novel strategy for solving optimization problem over the Stiefel manifold $St(p, N)$ via solving optimization problem over a vector space.
- To this end, we newly define a surjection $\mathfrak{C}: Q \rightarrow St(p, N)$ from a vector space Q onto $St(p, N)$, called a global Cayley parametrization.
- We present useful properties of \mathfrak{C} for optimization, e.g., the Lipschitz continuity of the gradient of $f \circ \mathfrak{C}$.
- Our numerical experiment verifies efficacies of the proposed strategy.

Optimization over $St(p, N) := \{U \in \mathbb{R}^{N \times p} \mid U^T U = I_p\}$

For a given continuous function $f: \mathbb{R}^{N \times p} \rightarrow \mathbb{R}$,

$$\text{find } U^* \in \operatorname{argmin}_{U \in St(p, N)} f(U) \dots (1)$$

Applications

Sparse PCA, Joint diagonalization, Enhancement of generalization in Deep Neural Network ...

Main difficulties

- $St(p, N)$ is not a vector space.
 - $\alpha U_1 + \beta U_2 \notin St(p, N)$ ($U_1, U_2 \in St(p, N), \alpha, \beta \in \mathbb{R}$)
- Many optimization techniques heavily rely on the linearity.
 - Gradient descent method: $x_{n+1} := x_n + \gamma_n d_n$ ($x_{n+1}, x_n, d_n \in \mathbb{R}^l, \gamma_n \in \mathbb{R}$)

Cayley-type transform Φ_S for $St(p, N)$ [Kume-Yamada'20, EUSIPCO]

Definition

For $S \in O(N) := St(N, N)$, let $E_{N,p}(S) := \{U \in St(p, N) \mid \det(I + S_{le}^T U) = 0\}$.

We define a mapping Φ_S as

$$\Phi_S: St(p, N) \setminus E_{N,p}(S) \rightarrow Q_{N,p} := \left\{ \begin{bmatrix} A & -B^T \\ B & 0 \end{bmatrix} \mid \begin{matrix} A^T = -A \in \mathbb{R}^{p \times p} \\ B \in \mathbb{R}^{(N-p) \times p} \end{matrix} \right\}$$

$$X \mapsto \begin{bmatrix} A_S(X) & -B_S^T(X) \\ B_S(X) & 0 \end{bmatrix} \quad \text{vector space}$$

$$S := \begin{bmatrix} S_{le} & S_{ri} \end{bmatrix} \quad N$$

p
 $N-p$

- $A_S(X) := (I + S_{le}^T X)^{-T} (X^T S_{le} - S_{le}^T X) (I + S_{le}^T X)^{-1} \in Q_{p,p}$
- $B_S(X) := -S_{ri}^T X (I + S_{le}^T X)^{-1} \in \mathbb{R}^{(N-p) \times p}$.

- Φ_S is diffeomorphic between the vector space $Q_{N,p}$ and $St(p, N) \setminus E_{N,p}$ with $\Phi_S^{-1}: Q_{N,p} \rightarrow St(p, N) \setminus E_{N,p}(S): V \mapsto S(I - V)(I + V)^{-1} I_{le}$.

- $\varphi := \Phi_I$ with $p = N$ is the classical Cayley transform for $SO(N) := \{U \in O(N) \mid \det(U) = 1\}$.

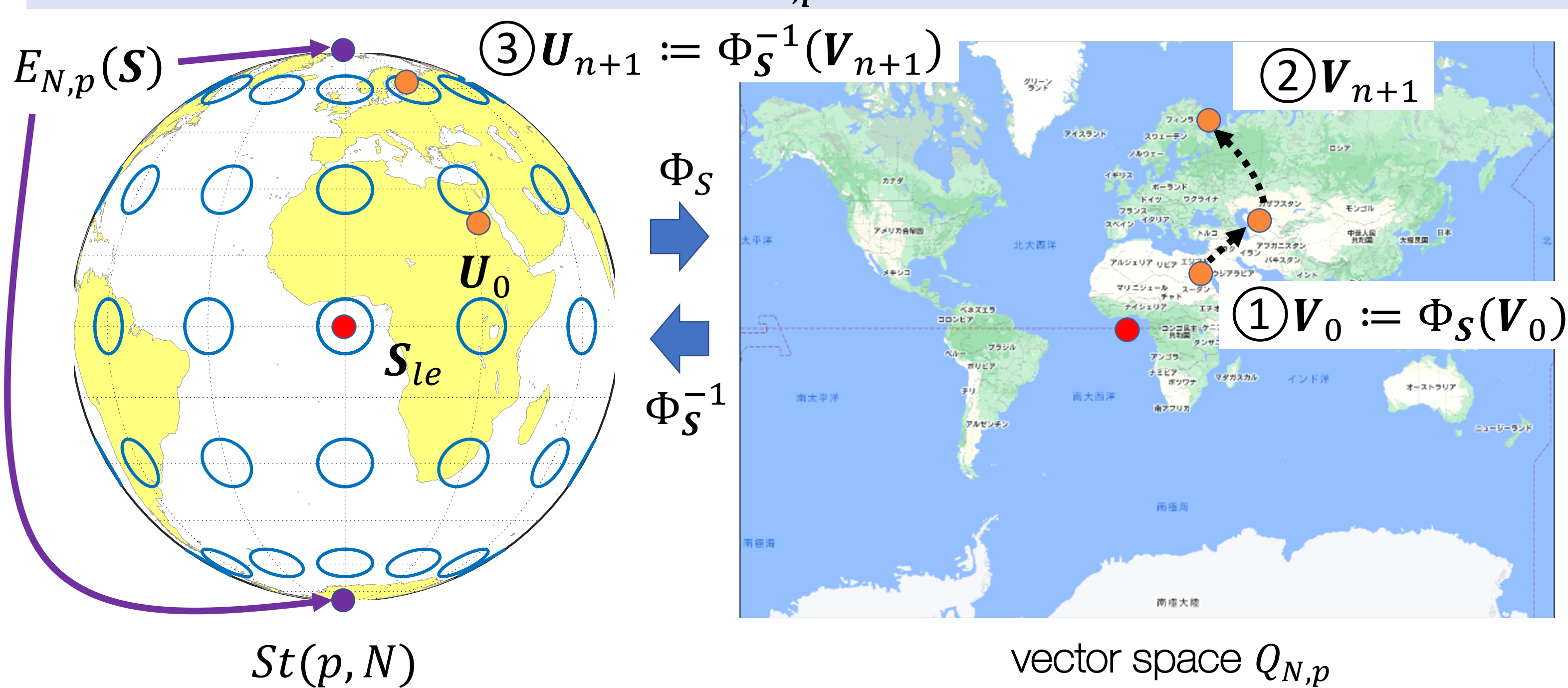
- S of Φ_S is called a center point because $\Phi_S(S_{le}) = 0$

Cayley parametrization (CP) strategy

Solve (1) via the following optimization problem over a vector space $Q_{N,p}$:

For a given continuous function $f: \mathbb{R}^{N \times p} \rightarrow \mathbb{R}$,

$$\text{find } V^* \in \operatorname{argmin}_{V \in Q_{N,p}} f \circ \Phi_S^{-1}(V) \dots (2)$$



- ① Translate an initial point $U_0 \in St(p, N)$ into $V_0 := \Phi_S(U_0) \in Q_{N,p}$.
- ② Update candidate solutions $(V_k)_{k=0}^{n+1} \subset Q_{N,p}$ in a vector space $Q_{N,p}$.
- ③ Translate the solution V_{n+1} into $U_{n+1} := \Phi_S^{-1}(V_{n+1}) \in St(p, N)$.

- We can utilize for ② directly optimization techniques over a vector space.

- $\Phi_S^{-1}(Q_{N,p}) \subsetneq St(p, N)$ for all $S \in O(N)$ possibly induces $\min f(St(p, N)) \neq \min f \circ \Phi_S^{-1}(Q_{N,p})$.

Natural question

Can we parameterize $St(p, N)$ by a single vector space entirely?

This study

Global Cayley parametrization (G-CP) strategy

Solve (1) via the following optimization problem over a vector space Q :

For a given continuous function $f: \mathbb{R}^{N \times p} \rightarrow \mathbb{R}$,

$$\text{find } V^* \in \min_{V \in Q} f \circ \mathfrak{C}(V) \dots (3)$$

Definition

Let $Q := Q_{p+1,p+1} \times Q_{p+1,p+1} \times Q_{N,p}$.

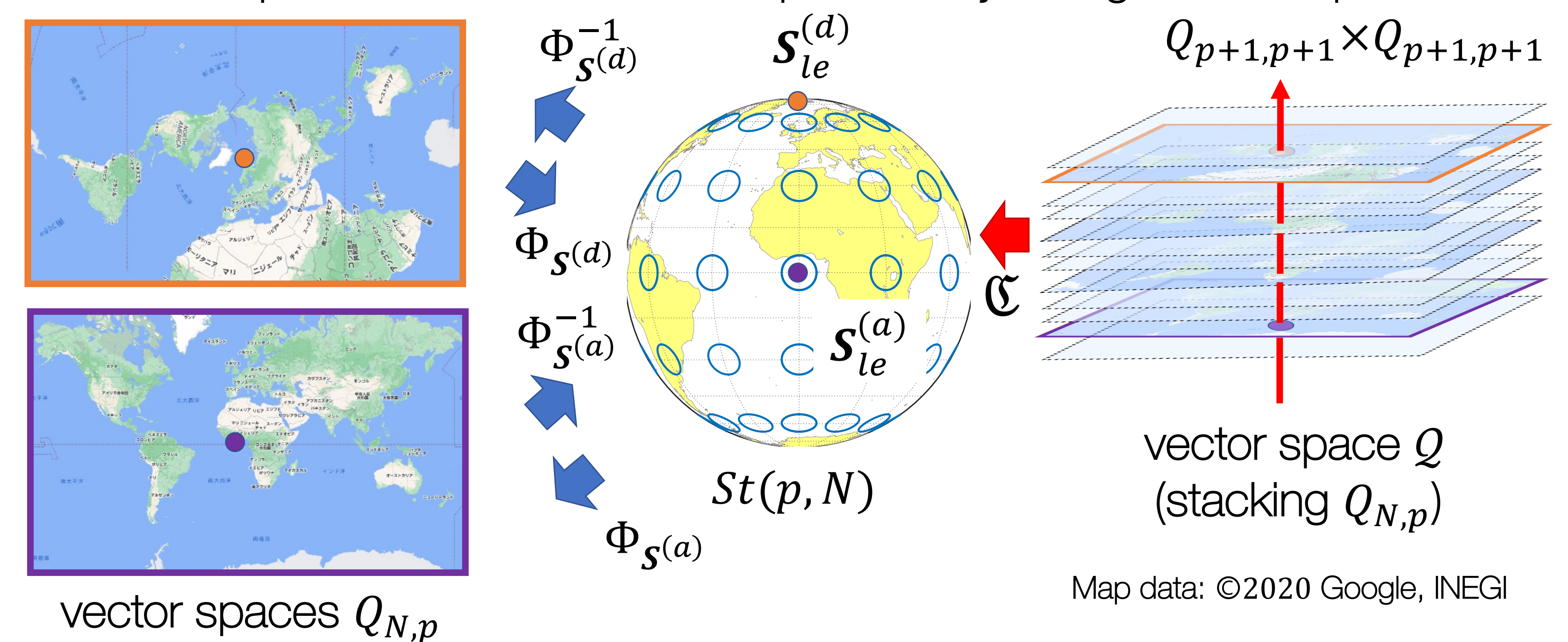
Then, the global Cayley parametrization is defined by

$$\mathfrak{C}: Q \rightarrow St(p, N): V := (V_1, V_2, V_3) \mapsto \Phi_{S(V_1, V_2)}^{-1}(V_3),$$

where $S(V_1, V_2) := \begin{bmatrix} \varphi^{-1}(V_1)\varphi^{-1}(V_2) & 0 \\ 0 & I_{N-p-1} \end{bmatrix} \in SO(N)$ and

$\varphi^{-1}: Q_{p+1,p+1} \rightarrow SO(p+1) \setminus E_{p+1,p+1}(I): V \mapsto (I - V)(I + V)^{-1}$ is the inversion mapping of the classical Cayley transform φ .

Basic idea: parameterization of center points S by a single vector space



vector spaces $Q_{N,p}$

Map data: ©2020 Google, INEGI

$$\text{Let } \widehat{SO}_{p+1}(N) := \left\{ \begin{bmatrix} \hat{T} & 0 \\ 0 & I_{N-p-1} \end{bmatrix} \in SO(N) \mid \hat{T} \in SO(p+1) \right\}.$$

It holds $\{\Phi_S^{-1}(V) \mid (S, V) \in \widehat{SO}_p(N) \times Q_{N,p} \cong SO(p+1) \times Q_{N,p}\} = St(p, N)$.

Combined with the following lemma, we can derive the surjection of \mathfrak{C} .

Lemma [Weyl'46, The Classical Groups]

$$\{\varphi^{-1}(V_1)\varphi^{-1}(V_2) \mid (V_1, V_2) \in Q_{p+1,p+1} \times Q_{p+1,p+1}\} = SO(p+1)$$

Characterization of local minimizer and stationary point by $f \circ \mathfrak{C}$

Key properties for \mathfrak{C} (Theorem 3.2)

1. For $U^* \in St(p, N)$, U^* is a local minimizer of (1) $\Leftrightarrow V^* \in Q$ s.t. $U^* = \mathfrak{C}(V^*)$ are local minimizers of (3).
2. Let f be differentiable. If $\nabla(f \circ \mathfrak{C})(V) = 0$ for $V \in Q$, then $\mathfrak{C}(V)$ is a stationary point of (1).

We can find a local minimizer (stationary point) of the problem (1) via the problem (3) defined over the vector space Q .

Useful properties of $f \circ \mathfrak{C}$ for optimization (Proposition 3.5)

Ex: Lipschitz continuity of $\nabla(f \circ \mathfrak{C})$

Let $f: \mathbb{R}^{N \times p} \rightarrow \mathbb{R}$ be continuously differentiable and $\max \|\nabla f(St(p, N))\|_F \leq \mu$.

Suppose ∇f is Lipschitz continuous with $L > 0$ over $St(p, N)$.

Then, $\nabla(f \circ \mathfrak{C})$ is Lipschitz continuous with $24(L + \mu)$ over Q .

Numerical experiments

- Optimization technique: gradient descent method (GDM)
- Comparison with
 1. CP strategy [Kume-Yamada'20]
 2. Retraction-based strategy with QR decomposition (QR) [Boumal et al.'14, J. Mach. Res.]

1. Joint diagonalization problem

$$f(U) := \sum_{i \in J} \|U^T A_i U - \text{Diag}(U^T A_i U)\|_F^2$$

- $N = 1000, p = 10, |J| = 10$

2. Eigenvalue problem

$$f(U) := -\text{Tr}(U^T A U)$$

- $N = 2000, p = 10$

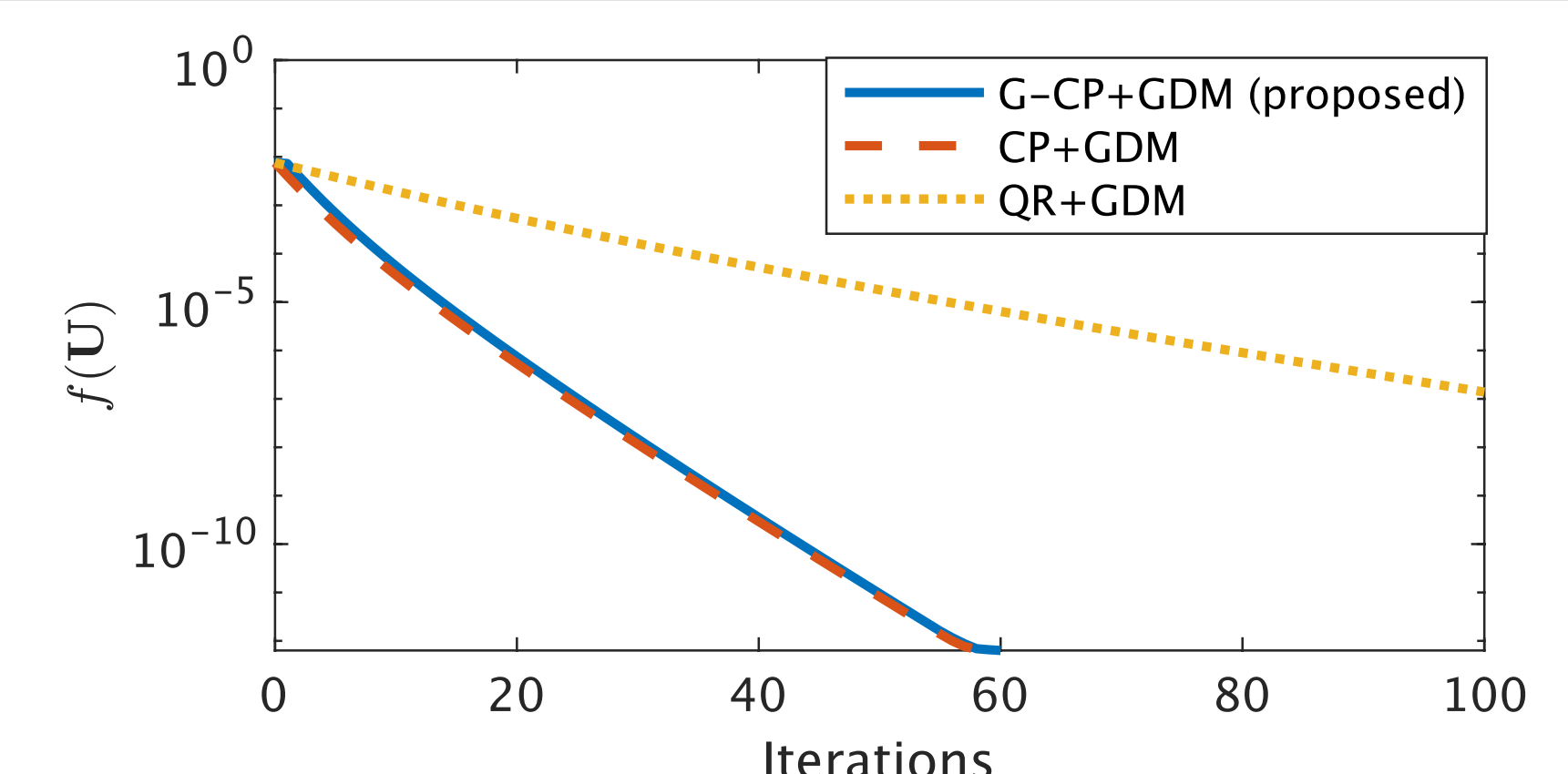


Fig.1 Result of the joint diagonalization problem

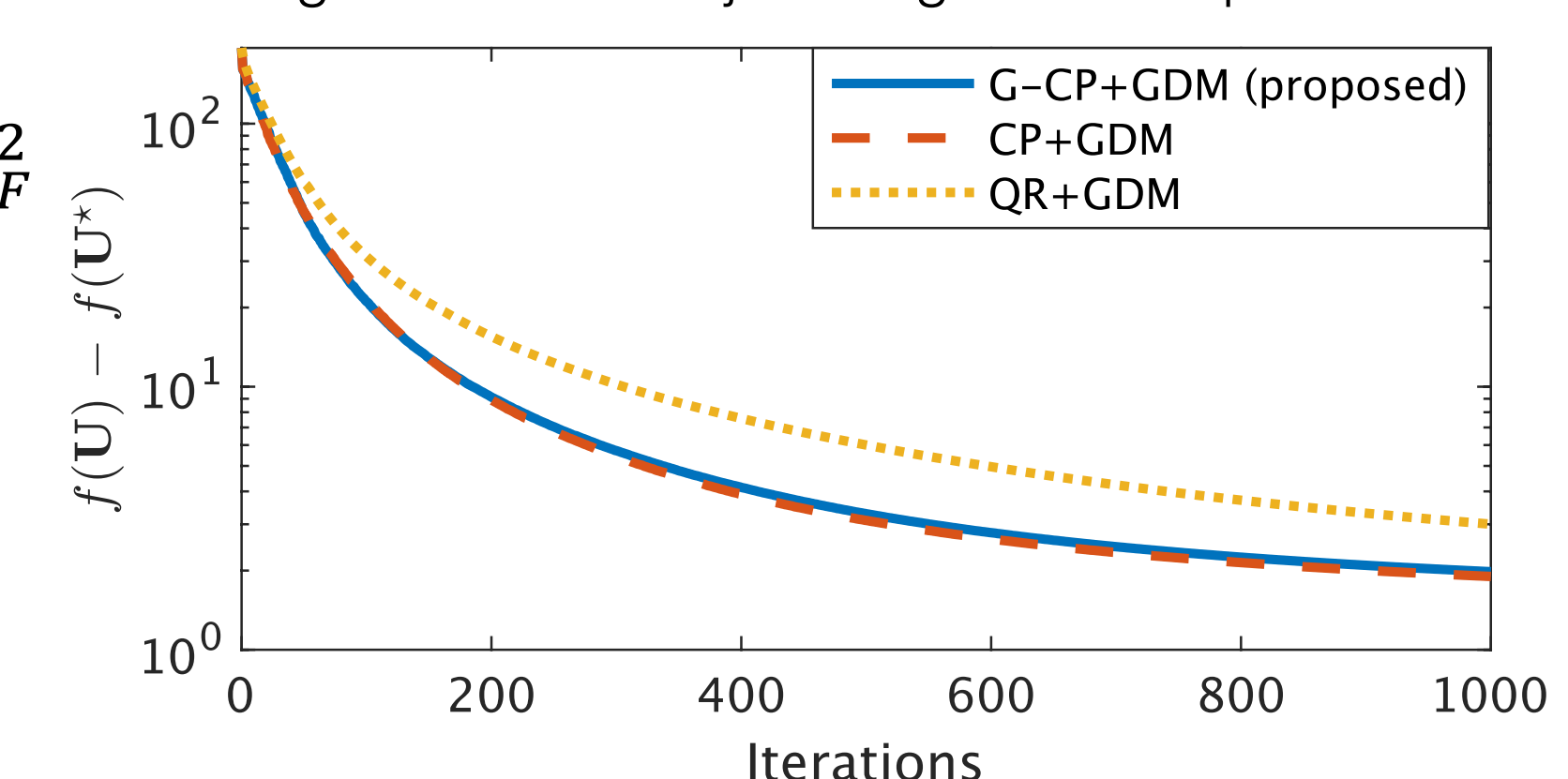


Fig.2 Result of the eigenvalue problem

G-CP strategy has potential to bring numerous optimization mechanisms over a vector space to (1) without losing the performance compared with CP strategy and the retraction-based strategy.