# OPTIMAL TRANSPORTATION FOR THE FAR-FIELD REFLECTOR PROBLEM

GANG BAO AND YIXUAN ZHANG

ABSTRACT. The inverse reflector problem aims to design a freeform reflecting surface that can direct the light from a specified source to produce the desired illumination in the target area, which is significant in the field of geometrical non-imaging optics. Mathematically, it can be formulated as an optimization problem, which is exactly the optimal transportation problem (OT) when the target is in the far field. The gradient of OT is governed by the generalized Monge-Ampère equation that models the far-field reflector system. Based on the gradient, this work presents a Sobolev gradient descent method implemented within a finite element framework to solve the corresponding OT. Convergence of the method is established and numerical examples are provided to demonstrate the effectiveness of the method.

## 1. INTRODUCTION

1.1. **Far-field reflector problem.** The far-field reflector system consists of a point source of the light placed at the origin $O$, a perfectly reflecting surface $\Gamma$ that is radial relative to $O$, and a far-field target area that receives the light originating from the source. The light source has the power density $f(x)$ in the direction $x \in \Omega \subset \mathbb{S}^2$, which describes the distribution of the intensity of rays from the source $O$. The reflecting surface $\Gamma$ can be characterized by

$$\Gamma := \{x\rho(x) : x \in \Omega, \ \rho > 0\}.$$

For a ray originating from $O$ in the direction $x$, it hits the surface $\Gamma$ at the point $x\rho(x)$ and produces a reflected ray in the direction $y \in \mathbb{S}^2$, which follows Snell's law

$$y = T(x) := x - 2(x \cdot \mathrm{n})\mathrm{n}, \tag{1.1}$$

where n is the outward normal of $\Gamma$ at the point $x\rho(x)$. Over all directions $x \in \Omega$, this reflecting procedure creates a far-field illumination intensity $g(y)$ on the domain $\Omega^* \subset \mathbb{S}^2$. Suppose there is no loss of energy in the reflection. Then by the energy conservation law, the map $T$ defined by (1.1) is measure preserving, i.e.,

$$\int_{T^{-1}(B)} f(x)\mathrm{d}x = \int_B g(y)\mathrm{d}y, \quad \forall \text{ Borel set } B \subset \Omega^*, \tag{1.2}$$

which is denoted by $T_\# f = g$, the pushforward measure of $f$. Suppose that the source intensity $f$ and the target intensity $g$ have been given in advance. For the above far-field reflector system, our goal is to reconstruct the reflecting surface $\Gamma$ that is capable of producing the prescribed illumination distribution $g(y)$ on the region $\Omega^*$ of a far-field sphere (see Figure 1).
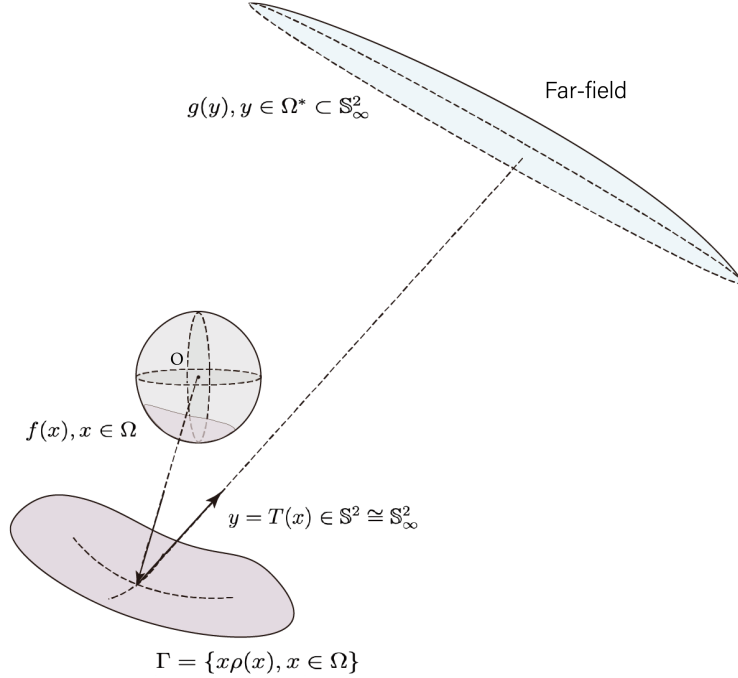
---

FIGURE 1. Far-field reflector system. The reflecting surface $\Gamma$ receives the source intensity $f(x)$ and produces the target intensity $g(y)$, where $x, y \in \mathbb{S}^2$.

The measure-preserving property (1.2) can be expressed as the Jacobian equation following by applying the change of variables,

$$g\left(T(x)\right)\left|\det(\nabla T)\right| = f(x). \tag{1.3}$$

The map $T$ depends on the radial distance function $\rho$ through the relation $\mathrm{n} = (\nabla\rho - \rho x)/\sqrt{\rho^2 + |\nabla\rho|^2}$. Hence by (1.1),

$$T(x) = \frac{\nabla\rho - (\eta - \rho)\,x}{\eta}, \qquad \eta = \left(|\nabla\rho|^2 + \rho^2\right)/2\rho. \tag{1.4}$$

Through (1.3), [37] yields a fully nonlinear partial differential equation in local coordinates,

$$\frac{1}{\eta^2}\det\left(-D^2\rho + \frac{2}{\rho}D\rho \otimes D\rho - (\rho - \eta)I\right) = \frac{f(x)}{g(T(x))}, \qquad x \in \Omega,$$
$$T(\Omega) = \Omega^*. \tag{1.5}$$

The existence, uniqueness and regularity of (1.5) have been studied in [11, 37, 38].

The numerical strategies for solving the far-field reflector problem can be roughly categorized as the ray-mapping method [6, 13, 16, 18], the supporting quadric method [12, 14, 22, 23, 30], and the optimal transport method. The ray-mapping method consists in building a ray mapping between the source and the target, and then constructing the surface that is capable of producing this map following (1.4). However, an integrability condition should be satisfied to ensure the existence of the surface, which is difficult to verify theoretically. For the supporting quadric method, the optical surface is constructed by taking the envelope of paraboloids of revolution that are controlled by points in the target area, leading to a computational complexity up to $O\left(\frac{N^4}{\epsilon}\log\frac{N^2}{\epsilon}\right)$ [23].

It was proved in [19, 38] that the far-field reflector problem is equivalent to the optimal transport problem with a logarithmic type cost function. In fact, this OT is a linear assignment problem, which can be solved using the auction algorithm [9] or the Hungarian algorithm with a complexity of $O(N^3)$. On the other hand, the entropic regularized OT can be solved by the Sinkhorn algorithm [4], which reduces the computational complexity to $O(N^2)$. However, the existence of the regularizer severely degrades the accuracy of the solution. Using OT, the equation (1.5) can be simplified into a generalized Monge-Ampère equation. In [34], the equation is solved by the least-squares method, which is stable but computationally complex and slow to converge. The work [20] offered a convergent finite difference method. However, due to the complexity of constructing a monotonic scheme in local coordinates and the slow convergence rate of the first-order Euler method, this approach results in a high computational cost. In [8] and [39], the equation is discretized using the B-spline scheme and the finite difference scheme, respectively, and solved with Newton-type methods. The convergence of these methods is highly dependent on the choice of the initial value and the regularity of the illumination distributions.

1.2. **Main results.** The existing methods focus on solving the standard Monge-Ampère equation $\det (I + D^2u) = f$ on the plane [3, 17, 24, 32, 35]. In [1], a first-order relaxation method has been proposed to solve the standard Monge-Ampère equation, given by

$$u^{k+1} = u^k + \Delta^{-1}\text{MA}\left(x, \nabla u^k, D^2u^k\right), \tag{1.6}$$

where $\text{MA}(x, \nabla u, D^2u) = f - \det (I + D^2u)$. In particular, (1.6) is implemented using the finite element method, and its convergence is established by demonstrating that (1.6) constitutes a fixed point scheme. Later, [29] extended this approach numerically to solve the $L^2$-Monge-Ampère equation on the sphere using the finite element method, addressing the challenge of mesh adaptivity. A similar approach was employed in [2] to compute the quadratic Wasserstein distance for applications in inverse problems. However, the generalized Monge-Ampère equation associated with (1.5) is significantly more complicated, where

$$\text{MA}(x, \nabla u, D^2u) = f(x) - \tilde{g}(x, \nabla u) \det \left(D^2u + A(x, \nabla u)\right). \tag{1.7}$$

The high nonlinearity of the equation makes the previous fixed point method ineffective for proving the convergence of (1.6) in the case of (1.7). Moreover, no convergence analysis is available, further complicating its numerical treatment.

In this work, we transform the far field reflector problem (1.5) into an optimal transport problem, which is equivalent to a generalized Monge-Ampère equation. We apply (1.6) to solve this equation and rigorously establish the convergence analysis of the method. By introducing the duality theory of OT, we demonstrate that the generalized Monge-Ampère operator in (1.7) corresponds to the gradient of the dual functional of OT. Consequently, $(-\Delta)^{-1}\text{MA}$ is identified as the $\dot{H}^1$ gradient in this context, and (1.6) is shown to be a descent scheme that guarantees the strict decrease of the dual functional, ultimately converging to the solution of the equation. Based on this, we refer to this approach as the Sobolev gradient method. It should be noted that our analysis also provides a novel convergence proof for the scheme (1.6) in the context of Monge-Ampère equations discussed in [1, 29]. Numerical experiments of this approach for the far-field problem can be implemented by adapting the finite element method from [29].

There are several advantages associated with the proposed method. First, (1.6) is a Poisson equation, which can be solved by numerous fast algorithms with a computational

complexity as low as $O(N)$. Both the theoretical proof and the numerical results can verify that the Sobolev gradient speeds up the convergence. Besides, different from the Newton-type methods, (1.6) does not need to compute the Jacobian matrix, which could be nontrivial for singular intensities. Compared with other methods, the approach proposed in this work is easier to accomplish, avoiding the complicated pre-calculations of numerical schemes. It provides an efficient way to solve the freeform reflector problem.

The paper is organized as follows. Section 2 introduces the relation between the optimal transport problem and the far-field reflector design. Section 3 derives the form of (1.7), which is related to the gradient of the dual functional of OT. Section 4 is devoted to proving the convergence of the Sobolev gradient descent scheme (1.6). In Section 5, the numerical algorithm for solving the far-field reflector problem is introduced. Section 6 provides some numerical examples to demonstrate the feasibility and effectiveness of the method.

**Notation and preliminaries.** For a matrix $\omega$, its subscripts $\omega_{ij}$ denotes the $(i,j)^{th}$ entry of $\omega$ and superscripts $\omega^{ij}$ denotes the $(i,j)^{th}$ entry of $\omega^{-1}$. For a matrix function $\omega(x) \in C(\Omega)^{n \times n}$, we define its norm $\|\omega\|_{p,\Omega}$ by

$$\|\omega\|_{p,\Omega} := \sup_{x \in \Omega} \|\omega(x)\|_p, \quad \text{where } \|\omega(x)\|_p := \sup_{b \in \mathbb{R}^n / 0} \frac{\|\omega(x) \cdot b\|_p}{\|b\|_p} \tag{1.8}$$

is induced by the $p$-norm for vectors.

The notation $D_i$ denotes the $i^{th}$ derivative in local coordinates, and the variable is specified in the subscript when the derivative operates on a multivariable function. In particular, for the cost function $c(x,y)$ in the optimal transport problem, we use the notation

$$c_{ij\cdots,kl\cdots} := D_{x_i} D_{x_j} \cdots D_{y_k} D_{y_l} \cdots c, \tag{1.9}$$

to simplify the notation of derivatives. In addition, the superscripts $c^{i,j}$ denotes the $(i,j)^{th}$ entry of the inverse matrix of $D_{xy}^2 c = (c_{i,j})$. The function space

$$\tilde{C}^{k,\alpha}(\Omega) := \left\{ u \in C^{k,\alpha}(\Omega) : \int_\Omega u(x)\mathrm{d}x = 0 \right\}, \tag{1.10}$$

is used to denote the $C^{k,\alpha}$ space with zero mass on the domain.

## 2. Freeform design and Optimal Transport

This section briefly reviews the basic concepts of the optimal transport problem and its connection with the far-field reflector design.

2.1. **Optimal transport problem.** Suppose that $\mathcal{X}$, $\mathcal{Y}$ are complete and separable metric spaces and $f$, $g$ are probability density functions defined on $\mathcal{X}$ and $\mathcal{Y}$, respectively. The cost function $c(x,y) : \mathcal{X} \times \mathcal{Y} \to \mathbb{R}$ indicates the cost of transporting a unit of mass from $x$ to $y$. Then Monge's optimal transport problem seeks the most efficient way of rearranging $f$ into $g$, which is to minimize

$$\inf_{T \in \Pi(f,g)} \int_{\mathcal{X} \times \mathcal{Y}} c(x, T(x)) f(x)\mathrm{d}x, \tag{2.1}$$

where $\quad \Pi(f, g) := \Big\{ T : \mathcal{X} \to \mathcal{Y} : \ T_{\#}f = g, \ \text{i.e.},$

$$\int_{T^{-1}(B)} f(x) = \int_B g(y), \quad \forall B \subset \mathcal{Y} \Big\},$$

is the set of all measure-preserving maps. The existence and uniqueness of the optimal map in (2.1) have been established for different cost functions on the Euclidean space and compact manifolds, including $c(x, y) = \frac{1}{2}|x - y|^2$ on $\mathbb{R}^d$ [7, 28], $c(x, y) = -\log(1 - x \cdot y)$ on $\mathbb{S}^2$, etc.

The dual formulation of the optimal transport problem is

$$\inf_{(u,v)} I(u, v) := \int_{\mathcal{X}} u(x)f(x)\mathrm{d}x + \int_{\mathcal{Y}} v(y)g(y)\mathrm{d}y,$$
$$(u, v) \in C(\mathcal{X}) \times C(\mathcal{Y}), \quad u(x) + v(y) + c(x, y) \geq 0. \tag{2.2}$$

This is a linear programming problem, which is desirable for designing numerical algorithms. The standard duality result shows that the infimum in (2.1) is equal to the negative of the infimum in (2.2).

2.2. **OT interpretation of the far-field reflector problem.** Considering $\mathcal{X} = \mathcal{Y} = \mathbb{S}^2$, the probability density functions $f$ and $g$ are the source intensity and target intensity in the far-field reflector problem, respectively. Then, the far-field problem (1.5) is equivalent to the optimal transport problem.

**Theorem 2.1** ( [19, 38]). *Suppose that $f$ and $g$ are bounded positive functions on the connected domains $\Omega \subset \mathbb{S}^2_-$ and $\Omega^* \subset \mathbb{S}^2_+$ respectively. Then there is a minimizer $(u, v)$ of the dual problem (2.2) for the cost function $c(x, y) = -\log(1 - x \cdot y)$. In particular, $(u, v)$ is unique up to a constant, and $\rho = e^{-u}$ solves the reflector problem (1.5).*

Therefore, by substituting $\rho = e^{-u}$ into (1.4), we can obtain the optical map in terms of the variable $u$,

$$T(x) = \frac{(|\nabla u|^2 - 1)x - 2\nabla u}{|\nabla u|^2 + 1}, \tag{2.3}$$

which depends only on the gradient of $u$.

In addition to the cost $-\log(1 - x \cdot y)$, the far-field problem can also be associated with (2.2) under the cost $c(x, y) = \log(1 - x \cdot y)$.

**Theorem 2.2** ( [38]). *Under the same assumptions of Theorem 2.1, there is a minimizer $(\varphi, \psi)$ of the dual problem (2.2) for $c(x, y) = \log(1 - x \cdot y)$. The minimizer is unique up to a constant and $\rho = e^{\varphi}$ solves the reflector problem (1.5).*

According to [38], the solution of (1.5) must be a constant multiplication of $\rho = e^{-u}$ of Theorem 2.1 or $\rho = e^{\varphi}$ of Theorem 2.2. In this work, we mainly focus on the case $\rho = e^{-u}$. Here, the minimizer $u$ is Lipschitz continuous but not necessarily $C^1$ smooth, and thus $\rho = e^{-u}$ is understood as a generalized solution [37] of (1.5). The regularity theorem [27] of OT implies that the smoothness of $u$ depends on the smoothness of $f, g$ and geometric structures of $\Omega, \Omega^*$.

The relationship between the far-field reflector problem and the primal form (2.1) of OT can also be developed, as detailed in Corollary 2.3. It is a direct result of OT's duality theory.

**Corollary 2.3** ( [19, 38]). *Suppose the same conditions as in Theorem 2.1 hold. Then the optical map $T$ given in (2.3) is the unique measure-preserving map minimizing the functional (2.1) for the cost function $c(x, y) = -\log(1 - x \cdot y)$.*

In this work, we aim to solve the freeform design problem for the far-field problem using a gradient descent method minimizing (2.2).

## 3. GRADIENT METHOD OF THE DUAL PROBLEM

According to Theorem 2.1, by substituting $\rho = e^{-u}$, the problem (1.5) is transformed into a generalized Monge-Ampère equation (3.1) in (1.7),

$$f(x) - \tilde{g}(x, \nabla u) \det \left( D^2 u + A(x, \nabla u) \right) = 0. \tag{3.1}$$

The explicit expressions of $\tilde{g}$ and $A$ are given Proposition 3.3 of this section.

Next, we prove that the expression of (3.1) is exactly the gradient of the dual functional and thus (3.1) can be viewed as the first order optimality condition of (2.2). Based on this fact, we can design gradient methods for optimizing (2.2).

Suppose that $\mathcal{X} = \mathcal{Y}$ is the sphere $\mathbb{S}^n$ or the flat torus $\mathbb{T}^n = \mathbb{R}^n/\mathbb{Z}^n$, $n \in \mathbb{Z}^+$. To ensure the regularity of the solution to the optimal transport problem, we assume that the cost function $c(x, y) \in C^{4 \times 4}(\mathcal{X} \times \mathcal{Y})$ is locally semi-concave, and satisfies Ma-Trudinger-Wang (MTW) condition in [27, 36] as well as the twist condition,

- $x \mapsto \nabla_x c(x, y)$ is injective for $(x, y) \in \mathcal{X} \times \mathcal{Y}$;
- $\det D^2_{xy} c(x, y) \neq 0$ on $\mathcal{X} \times \mathcal{Y}$.

First, we introduce the notation of $c$-transform, which provides an explicit relationship between $u$ and $v$ in (2.2), thereby defining the transport map and simplifying the optimization form of (2.2).

**Definition 3.1.** For $u : \mathcal{X} \to \mathbb{R}$ and $v : \mathcal{Y} \to \mathbb{R}$, their $c$-transforms are defined by

$$\begin{aligned} u^c(y) &= \sup_{x \in \mathcal{X}} -c(x, y) - u(x) \\ v^c(x) &= \sup_{y \in \mathcal{Y}} -c(x, y) - v(y) \end{aligned} \tag{3.2}$$

A function $u$ (or $v$) is said to be $c$-convex if $u = u^{cc}$ (or $v = v^{cc}$).

It can be verified that $c$-convex functions are Lipschitz continuous [36]. For a $c$-convex function $u$ and any $x_0 \in \mathcal{X}$, there exists $y_0 \in \mathcal{Y}$ such that

$$\begin{cases} u(x_0) + u^c(y_0) + c(x_0, y_0) &= 0, \\ u(x) + u^c(y_0) + c(x, y_0) &\geq 0, \quad \forall x \in \mathcal{X}. \end{cases} \tag{3.3}$$

If $u \in C^2(\mathcal{X})$, the optimality conditions of (3.3) imply that

$$\nabla u(x_0) + \nabla_x c(x_0, y_0) = 0, \tag{3.4}$$

$$D^2 u(x_0) + D^2_{xx} c(x_0, y_0) \geq 0. \tag{3.5}$$

From (3.4), we can define a map $T_u : \mathcal{X} \to \mathcal{Y}$ by

$$T_u(x) = \mathcal{T}(x, \nabla u), \quad \text{where } \mathcal{T}(x, p) := (\nabla_x c)^{-1}(x, -p), \tag{3.6}$$

which maps each $x_0$ to the corresponding $y_0$. In fact, based on (3.5), we can also propose a sufficient condition to ensure the $c$-convexity of $u \in C^2(\mathcal{X})$, stated in the following lemma.

**Lemma 3.2** ( [21]). *Suppose that $u \in C^2(\mathcal{X})$ and*

$$D^2 u(x) + A\left(x, \nabla u(x)\right) > \lambda, \quad \forall x \in \mathcal{X}, \tag{3.7}$$

*for some $\lambda > 0$, where $A(x, p) := D^2_{xx} c(x, \mathcal{T}(x, p))$. Then $u$ is c-convex and $T_u$ in (3.6) is a diffeomorphism. The equality $u(x) + u^c(y) + c(x, y) = 0$ holds if and only if $y = T_u(x)$.*

In general, $u \in C^2(\mathcal{X})$ is said to be strictly $c$-convex with respect to $\lambda$ if (3.7) is satisfied. For any function $u$, since $I(u^{cc}, u^c) \le I(u, u^c) \le I(u, v)$ , it is possible to reduce (2.2) to the optimization problem depending on a single variable $u$, see [36].

**Proposition 3.3.** *The dual problem (2.2) is equivalent to maximizing the functional $\mathcal{J}$ over the set of c-convex functions:*

$$\mathcal{J}(u) := \int_{\mathcal{X}} u(x) f(x) \mathrm{d}x + \int_{\mathcal{Y}} u^c(y) g(y) \mathrm{d}y. \tag{3.8}$$

*The functional $\mathcal{J}$ is convex and Lipschitz continuous. Moreover, if $u \in C^2(\mathcal{X})$ is strictly c-convex, the first variation of $\mathcal{J}$ at $u$ is given by*

$$\mathcal{J}'(u) = f - g(T_u) \left|\det(\nabla T_u)\right|, \tag{3.9}$$

*where $T_u = \mathcal{T}(x, \nabla u)$ is the map defined in (3.6). The formula (3.9) is equal to*

$$\mathcal{J}'(u) = f - \tilde{g}\left(x, \nabla u\right) \det\left(D^2 u + A(x, \nabla u)\right), \tag{3.10}$$

*where $A(x, p)$ is given in (3.6) and*

$$\tilde{g}(x, p) := \frac{g(\mathcal{T}(x, p))}{\left|\det\left(D^2_{xy} c\left(x, \mathcal{T}(x, p)\right)\right)\right|}.$$

If the $c$-convex function $u \in C^2(\mathcal{X})$ achieves the infimum of the functional $\mathcal{J}$, then $T_u$ is exactly the optimal transport map solving (2.1) and the generalized Monge-Ampère equation (3.1) holds. For $c(x, y) = -\log(1 - x \cdot y)$, it can be checked that the optimal map $T_u$ computed through (3.6) coincides the optical map (2.3).

*Proof of Proposition 3.3.* The convexity of the functional follows from

$$(tu_1 + (1 - t)u_2)^c(y) = \sup_{x \in \mathcal{X}} \left(-c(x, y) - tu_1(x) - (1 - t)u_2(x)\right)$$

$$\le \sup_{x \in \mathcal{X}} t\left(-c(x, y) - u_1(x)\right) + \sup_{x \in \mathcal{X}}(1 - t)\left(-c(x, y) - u_2(x)\right)$$

$$= tu_1^c(y) + (1 - t)u_2^c(y).$$

For any $u_1, u_2 \in C(\mathcal{X})$,

$$u_1^c(y) = \sup_x -c(x, y) - u_1(x) \ge \sup_x -c(x, y) - u_2(x) - \|u_2 - u_1\|_{C^0}$$

$$= u_2^c(y) - \|u_2 - u_1\|_{C^0},$$

which means that

$$\|u_1^c - u_2^c\|_{C^0} \le \|u_2 - u_1\|_{C^0}.$$

Further, we obtain the Lipschitz continuity of $\mathcal{J}$ on $C(\mathcal{X})$,

$$\left|\mathcal{J}\left(u_1\right) - \mathcal{J}\left(u_2\right)\right| \le \|f\|_{L^1} \|u_1 - u_2\|_{C^0} + \|g\|_{L^1} \|u_1^c - u_2^c\|_{C^0} \le 2 \|u_1 - u_2\|_{C^0}.$$

Considering the strictly $c$-convex function $u \in C^2(\mathcal{X})$. Let $v \in C(\mathcal{X})$ and $\epsilon$ be a small enough constant. For fixed $y \in \mathcal{Y}$, we introduce the notation $x := T_u^{-1}(y)$ and

$$x_\epsilon \in \mathrm{argmax}_{x \in \mathcal{X}} \left\{-c(x, y) - (u + \epsilon v)(x)\right\}.$$

Using the property of the $c$-transform, we obtain

$$(u + \epsilon v)^c(y) - u^c(y) = -c(x_\epsilon, y) - (u + \epsilon v)(x_\epsilon) - u^c(y)$$
$$\leq u(x_\epsilon) - (u + \epsilon v)(x_\epsilon)$$
$$= -\epsilon v(x_\epsilon),$$
$$(u + \epsilon v)^c(y) - u^c(y) = (u + \epsilon v)^c(y) + c(x_0, y) + u(x_0)$$
$$\geq -(u + \epsilon v)(x_0) + u(x_0)$$
$$= -\epsilon v(x_0).$$

It follows that

$$\left| \frac{(u + \epsilon v)^c(y) - u^c(y)}{\epsilon} + v(x) \right| \leq |v(x_\epsilon) - v(x)|.$$

Every convergent subsequence of $\{x_\epsilon\}$ converges to the point $x$, thus $\{x_\epsilon\} \to x$ as $\epsilon \to 0$. Therefore, $v(x_\epsilon) \to v(x)$ and

$$\lim_{\epsilon \to 0} \frac{(u + \epsilon v)^c(y) - u^c(y)}{\epsilon} = -(v \circ T_u^{-1})(y).$$

Hence we have

$$\mathcal{J}'(u)v = \lim_{\epsilon \to 0} \frac{\mathcal{J}(u + \epsilon v) - \mathcal{J}(u)}{\epsilon} = \int_{\mathcal{X}} v(x) f(x) \mathrm{d}x + \int_{\mathcal{Y}} \lim_{\epsilon \to 0} \frac{(u + \epsilon v)^c(y) - v^c(y)}{\epsilon} g(y) \mathrm{d}y$$

$$= \int_{\mathcal{X}} v(x) f(x) \mathrm{d}x - \int_{\mathcal{Y}} (v \circ T_u^{-1})(y) g(y) \mathrm{d}y$$

$$= \int_{\mathcal{X}} v(x) \big( f - g(T_u) | \det(\nabla T_u)| \big) \mathrm{d}x.$$

Thus we obtain $\mathcal{J}'(u)$ in (3.9). A a direct computation [26, 36] of (3.9) yields the formula (3.10).

$$\square$$

Furthermore, the operator $\mathcal{J}'(u)$ is differentiable on the space $C^2(\mathcal{X})$, which gives the second order information of $\mathcal{J}$.

**Corollary 3.4.** *The second order derivative of the functional $\mathcal{J}$ is given by the linearized operator $\mathcal{L}$ of $\mathcal{J}'$ at $u$,*

$$\mathcal{L}_u v := \sum_i \mathcal{L}_u^i D_i v + \sum_{ij} \mathcal{L}_u^{ij} D_{ij}^2 v, \quad v \in C^2(\mathcal{X}), \tag{3.11}$$

$$\mathcal{L}_u^i := \det(\omega) \sum_{kl} \left( D_{p_i} \tilde{g} + \tilde{g} \omega^{kl} D_{p_i} A_{kl} \right), \quad \mathcal{L}_u^{ij} := \tilde{g} \det(\omega) \omega^{ij},$$

*where $\omega := D^2 u + A(x, \nabla u)$ and $D_{p_i} A_{kl}, D_{p_i} \tilde{g}$ are evaluated at the point $(x, p) = (x, \nabla u)$. Here,*

$$D_{p_i} A_{kl}(x, p) = \sum_j c_{kl,j} c^{j,i},$$

$$\tag{3.12}$$

$$D_{p_i} \tilde{g}(x, p) = \left| \det(D_{xy}^2 c) \right|^{-1} \left( \sum_j c^{j,i} (D_j g \circ \mathcal{T}) + (g \circ \mathcal{T}) \sum_{jkl} c^{j,k} c_{k,jl} c^{l,i} \right).$$

*where the cost function $c$ and its derivatives are evaluated at $(x, y) = (x, \mathcal{T}(x, p))$.*

For $c(x,y) = \frac{1}{2}|x-y|^2$ on $\mathbb{R}^d$, the equation (3.1) is $g(x+\nabla u)\det(D^2u+I) = f(x)$, which can be solved by Newton's method using (3.11). However, for the far-field reflector problem $c(x,y) = -\log(1 - x \cdot y)$ on $\mathbb{S}^2$, due to the high nonlinearity of the equation and the domain, $\mathcal{L}$ easily becomes nontrivial. The gradient method is preferred for this case, i.e.,

$$u^{k+1} = u^k - \tau \cdot \mathcal{J}'(u^k).$$

However, the step size $\tau$ is hard to choose to ensure the convergence of the method. Motivated by the Sobolev gradient [31], we can define a new gradient $\mathcal{J}'_s$ by changing the underlying inner space from $L_2$ to $\dot{H}^1$ (the seminorm of $H^1$),

$$\langle \mathcal{J}', \eta \rangle_{L^2} = \langle \mathcal{J}'_s, \eta \rangle_{\dot{H}^1} := \langle \nabla \mathcal{J}'_s, \nabla \eta \rangle_{L^2}.$$

More formally, we define

$$\mathcal{J}'_s := (-\Delta)^{-1} \mathcal{J}',$$

which yields the Sobolev gradient descent scheme of (3.8),

$$\Delta u^{k+1} = \Delta u^k + \tau \cdot \mathcal{J}'(u^k). \tag{3.13}$$

In this work, we utilize the descent scheme (3.13) to develop the numerical algorithm for solving the far-field problem.

## 4. Convergence of the Sobolev descent scheme

This section is devoted to the convergence analysis of the Sobolev gradient descent sequence (3.13) for the optimization problem (3.8).

We assume that the conditions stated at the beginning of Section 3 are satisfied. We begin with the following regularity theorem [10, 21, 25–27, 36] of the generalized Monge-Ampère equation (3.1).

**Theorem 4.1.** *Suppose that $f, g \in C^{1,1}(\mathcal{X})$ have positive lower bounds on $\mathcal{X}$. Then the generalized Monge-Ampère equation (3.1) has a unique c-convex solution $u \in \tilde{C}^{3,\alpha}(\mathcal{X})$. In addition, $u$ is the unique c-convex minimizer of $\mathcal{J}(u)$ in (3.8).*

For this solution $u$, (3.5) implies that $\omega = D^2u + A(x, \nabla u)$ is a positive semi-definite matrix. Considering the boundedness of $\|u\|_{C^2}$, there exists a constant $C > 0$ such that

$$\|\omega\|_{2,\mathcal{X}} < C. \tag{4.1}$$

Since $c \in C^4$ and $\det(D^2_{xy}c) \neq 0$, we know that $\left|\det D^2_{xy}c\right|$ has a positive lower bound on $\mathcal{X} \times \mathcal{Y}$. Therefore, the boundedness of $f, g$ in Theorem 4.1 indicates that

$$\det(\omega) = \frac{f(x)}{\tilde{g}(x, \nabla u(x))} > C, \quad \forall x \in \mathcal{X}, \tag{4.2}$$

for some $C > 0$. The inequality (4.1) and (4.2) together imply that the eigenvalues of $\omega$ are bounded from above and below, i.e., there exists some $\lambda > 0$ such that

$$\frac{2}{\lambda} < D^2u(x) + A(x, \nabla u(x)) < \frac{1}{2}\lambda, \quad \forall x \in \mathcal{X}. \tag{4.3}$$

Obviously, there exists a neighborhood $V_\lambda$ of $u$ with the radius $r$

$$V_\lambda := \left\{ v \in \tilde{C}^2(\mathcal{X}) : \|u - v\|_{C^2(\mathcal{X})} \leq r(\lambda) \right\}, \tag{4.4}$$

such that a function $v$ in $V_\lambda$ is strictly $c$-convex with respect to $\lambda^{-1}$,

$$\frac{1}{\lambda} < D^2 v(x) + A(x, \nabla v(x)) < \lambda, \quad \forall x \in \mathcal{X}.$$

If the descent sequence is contained in this neighborhood of the solution $u$, then we can prove the strictly decreasing of $\mathcal{J}$ and the local convergence of $\{u^k\}_k$. To prove that, we need the following result to bound the divergence of the cofactor matrix of $\omega$ by the parameter $\lambda$ in (4.3).

**Lemma 4.2.** *Suppose that $u$ is an arbitrary strictly $c$-convex function. For the map $T_u$ defined in (3.6), its Jacobi matrix has divergence-free rows, i.e.*

$$\sum_{j=1}^n D_j \left( (\operatorname{cof} \nabla T_u)_{ij} \right) = 0, \qquad i = 1, \cdots, n. \tag{4.5}$$

*where $\operatorname{cof} \nabla T_u$ denotes the cofactor matrix of $\nabla T_u$. In addition, for any $x \in \mathcal{X}$, the matrix $\omega = D^2 u(x) + A(x, \nabla u(x))$ satisfies*

$$\left| \sum_{j=1}^n D_j \left( (\operatorname{cof} \omega)_{ij} \right) \right| \le C \left\| \operatorname{cof} \omega \right\|_{2,\mathcal{X}} \cdot \left( \left\| \omega \right\|_{2,\mathcal{X}} + 1 \right), \qquad i = 1, \cdots, n, \tag{4.6}$$

*where the constant $C$ depends only on the cost function $c$ and the dimension $n$.*

*Proof.* We refer to [15] for the proof of the identity (4.5). Here, we prove the inequality (4.6). It can be computed from (3.6) that

$$\nabla T_u(x) = -\left( D^2_{xy} c(x, T_u(x)) \right)^{-1} \left( D^2 u(x) + D^2_{xx} c(x, T_u(x)) \right). \tag{4.7}$$

For the cofactor matrix of a matrix $P$, we recall the identity $(\det P) I = P^{\mathrm{T}} (\operatorname{cof} P)$. Therefore, we obtain from (4.7)

$$\operatorname{cof} \nabla T_u = -\operatorname{cof} C \cdot \operatorname{cof} \omega,$$

where $\omega = D^2 u(x) + A(x, \nabla u(x))$ and $C := D^2_{xy} c(x, T_u(x))$. Hence,

$$-D_j \left( (\operatorname{cof} \nabla T_u)_{ij} \right) = \sum_{k=1}^n (\operatorname{cof} \omega)_{kj} \cdot D_j \left( (\operatorname{cof} C)_{ik} \right) + (\operatorname{cof} C)_{ik} \cdot D_j \left( (\operatorname{cof} \omega)_{kj} \right).$$

Summing the above formula over $j$, we have the equation

$$0 = \sum_{k=0}^n (\operatorname{cof} C)_{ik} \left( \sum_{j=1}^n D_j \left( (\operatorname{cof} \omega)_{kj} \right) \right) + \sum_{k,j=1}^n (\operatorname{cof} \omega)_{kj} \cdot D_j \left( (\operatorname{cof} C)_{ik} \right),$$

for $i = 1, 2, \cdots, n$. It can be reformulated as a linear equation

$$\mathbf{A}\mathbf{x} = \mathbf{b},$$

$$\mathbf{x} = (\mathrm{x}_1, \cdots \mathrm{x}_n) \text{ with } \mathrm{x}_i = \sum_{j=1}^n D_j \left( (\operatorname{cof} \omega)_{ij} \right), \qquad \mathbf{A} = (C^{\mathrm{T}})^{-1},$$

$$\mathbf{b} = (\mathrm{b}_1, \cdots \mathrm{b}_n) \text{ with } \mathrm{b}_i = -\frac{1}{\det C} \sum_{k,j=1}^n (\operatorname{cof} \omega)_{kj} \cdot D_j \left( (\operatorname{cof} C)_{ik} \right).$$

Here, the term $D_j \left( (\operatorname{cof} C)_{ik} \right)$ in $\mathbf{b}$ can be computed by the Jacobi's formula

$$D_j (\operatorname{cof} C) = (C^{\mathrm{T}})^{-1} \left( \operatorname{tr} \left( \operatorname{cof} C^{\mathrm{T}} \cdot D_j C \right) - D_j C^{\mathrm{T}} \cdot \operatorname{cof} C \right).$$

Through a simple computation, we have

$$
\mathrm{b}_i = -\sum_{kj}(\operatorname{cof}\omega)_{kj}\sum_{lmrs}\omega_{sj}c^{r,s}\big(c^{l,i}c^{k,m}c_{m,lr}-
$$
$$
c^{k,i}c^{m,l}c_{l,mr}\big) + \big(c^{k,i}c^{m,l}c_{lj,m} - c^{l,i}c^{k,m}c_{mj,l}\big)\,.
$$
(4.8)

The definitions of the notations have been given in (1.8) and (1.9). For the elements $\omega_{sj}$ and $(\operatorname{cof}\omega)_{kj}$ in (4.8),

$$
\sup_{x\in\mathcal{X}}\sup_{s,j}|\omega_{sj}| \le \|\omega\|_{2,\mathcal{X}}\,,
$$
$$
\sup_{x\in\mathcal{X}}\sup_{k,j}|(\operatorname{cof}\omega)_{kj}| \le \|\operatorname{cof}\omega\|_{2,\mathcal{X}}\,.
$$

Finally, we have the estimate

$$
\|\mathbf{x}\|_{\infty,\mathcal{X}} \le \left\|\mathbf{A}^{-1}\right\|_{\infty,\mathcal{X}}\|\mathbf{b}\|_{\infty,\mathcal{X}} \le C\,\|\operatorname{cof}\omega\|_{2,\mathcal{X}}\left(\|\omega\|_{2,\mathcal{X}}+1\right).
$$

Here, the constant $C$ depends only on the dimension $n$, the cost function $c$ up to its third-order derivatives and its twist condition,

$$
\|c\|_{C^{1\times2}(\mathcal{X}\times\mathcal{Y})}\,,\quad \|c\|_{C^{2\times1}(\mathcal{X}\times\mathcal{Y})}\,,\quad \left\|(D_{xy}c)^{-1}\right\|_{2,\mathcal{X}\times\mathcal{Y}}\,.
$$

Thus, we complete the proof of the inequality (4.6). $\qquad\square$

Now we are prepared to evaluate the decreasing of $\mathcal{J}$ on the sequence $\{u^k\}_k$ in (3.13).

**Proposition 4.3.** *Suppose that the same conditions as in Theorem 4.1 hold, and for sufficiently small $\tau > 0$, the sequence $\{u^k\}_k$ computed by*

$$
\Delta u^{k+1} = \Delta u^k + \tau\cdot\mathcal{J}'(u^k),
$$
(4.9)

*is uniformly contained in the neighborhood $V_\lambda$ defined by (4.4). Then there exists $\tau_{max} > 0$ such that when $\tau < \tau_{max}$, the functional $\mathcal{J}$ is strictly decreasing on $\{u^k\}_k$ and*

$$
\mathcal{J}(u^k) - \mathcal{J}(u^{k+1}) \ge \frac{1}{2\tau_{max}}\left\|\nabla u^k - \nabla u^{k+1}\right\|_{L^2}^2.
$$
(4.10)

*Here, $\tau_{max}$ depends on the dimension $n$, the parameter $\lambda$ in (4.3), the cost function $c$ and the density function $g$.*

*Proof.* From Proposition 3.3, the functional $\mathcal{J}(u)$ is convex and hence for the iteration (4.9), we have

$$
\mathcal{J}(u^k) - \mathcal{J}(u^{k+1}) \ge \int_{\mathcal{X}}\mathcal{J}'(u^{k+1})(u^k - u^{k+1})
$$
$$
= \frac{1}{\tau}\int_{\mathcal{X}}(u^k - u^{k+1})(\Delta u^{k+1} - \Delta u^k) + \int_{\mathcal{X}}(u^k - u^{k+1})\left(\mathcal{J}'(u^{k+1}) - \mathcal{J}'(u^k)\right).
$$
(4.11)

In order to estimate (4.11), we first consider

$$
\int_{\mathcal{X}}v\cdot\mathcal{L}_u v\,\mathrm{d}x = \sum_i\int_{\mathcal{X}}v\,\mathcal{L}_u^i\,D_i v + \sum_{ij}\int_{\mathcal{X}}v\,\mathcal{L}_u^{ij}\,D_{ij}^2 v,
$$
(4.12)

for any $u \in V_\lambda$ and $v \in \tilde{C}^2(\mathcal{X})$, where $\mathcal{L}_u$ is the linear operator given in Corollary 3.4. Since $\mathcal{X}$ is a closed manifold without boundary, applying the divergence theorem to the second term in (4.12), we obtain

$$\sum_{ij} \int_\mathcal{X} v\, \mathcal{L}_u^{ij}\, D_{ij}^2 v = \sum_{ij} \int_\mathcal{X} D_j v\, \mathcal{L}_u^{ij}\, D_i v + \sum_i \int_\mathcal{X} v \left( \operatorname{div} \mathcal{L}_u^{ij} \right) D_i v, \tag{4.13}$$

where

$$\operatorname{div} \mathcal{L}_u^{ij} := \sum_j D_j(\mathcal{L}_u^{ij}) = \sum_j \tilde{g}(x, \nabla u) D_j \left( (\operatorname{cof} \omega)_{ij} \right) + (\operatorname{cof} \omega)_{ij}\, D_j \left( \tilde{g}\,(x, \nabla u) \right) \tag{4.14}$$

and

$$D_j \left( \tilde{g}\,(x, \nabla u) \right) = \sum_{klrs} D_k g\, c^{k,l} \omega_{lj} - g\, c^{s,r} c_{rj,s} + g\, c^{s,r} c_{r,sk} c^{k,l} \omega_{lj}.$$

Using (4.6) in Lemma 4.2, the divergence (4.14) can be bounded by

$$|\operatorname{div} \mathcal{L}_u^{ij}| \leq C\|g\|_{C^1}\, \|\operatorname{cof} \omega\|_{2,\mathcal{X}} \left( \|\omega\|_{2,\mathcal{X}} + 1 \right),$$

where the constant $C$ depends only on the dimension $n$ and the cost function $c$ as in Lemma 4.2.

Similarly, we have $|\mathcal{L}_u^i| \leq C\|g\|_{C^1}$ from (3.12). Since $u \in V_\lambda$, we obtain $\|\omega\|_{2,\mathcal{X}} \leq \lambda$ and $\|\operatorname{cof} \omega\|_{2,\mathcal{X}} \leq \lambda^{n-1}$. Therefore, there exists an upper bound $C$ such that

$$\left| \mathcal{L}_u^i + \operatorname{div} \mathcal{L}_u^{ij} \right| \leq C(n, \lambda, c, g), \tag{4.15}$$

for any $x \in \mathcal{X}$ and $i, j = 1 \cdots, n$. The constant $C$ depends on the dimension $n$, the parameter $\lambda$, the cost function $c$ and the density function $g$.

Hence, combing with (4.13), we can estimate (4.12) with

$$\begin{aligned}
\left| \int_\mathcal{X} v \cdot \mathcal{L}_u v \, \mathrm{d}x \right| &= \left| \sum_i \int_\mathcal{X} v \left( \mathcal{L}_u^i + \operatorname{div} \mathcal{L}_u^{ij} \right) D_i v + \sum_{ij} \int_\mathcal{X} D_i v\, \mathcal{L}_u^{ij}\, D_j v \right| \\
&\leq C \sum_i \left( \int_\mathcal{X} |v| \cdot |D_i v| + \int_\mathcal{X} |D_i v|^2 \right) \\
&\leq C \left( \|v\|_{L^2} \sum_i \left( \int_\mathcal{X} |D_i v|^2 \right)^{\frac{1}{2}} + \|\nabla v\|_{L^2}^2 \right) \\
&\leq C \left( \|v\|_{L^2} \|\nabla v\|_{L^2} + \|\nabla v\|_{L^2}^2 \right) \\
&\leq C \|\nabla v\|_{L^2}^2,
\end{aligned} \tag{4.16}$$

where the constant $C$ depends on $n, \lambda, c, g$ as in (4.15). Here, the first inequality in (4.16) follows from (4.15) and the positive definiteness of the matrix $\mathcal{L}_u^{ij}$, which is bounded from above since $u \in V_\lambda$. The last inequality in (4.16) is based on the Poincaré–Wirtinger inequality for $v \in \tilde{C}^2(\mathcal{X})$.

For $u^k, u^{k+1} \in V_\lambda$, there exists $s \in (0, 1)$ such that

$$\mathcal{J}'(u^{k+1}) - \mathcal{J}'(u^k) = \mathcal{L}_u v, \tag{4.17}$$

where $v = u^{k+1} - u^k \in \tilde{C}^2$ and $u = s u^k + (1 - s) u^{k+1} \in V_\lambda$.

We can substitute (4.17) into (4.16) and obtain

$$\left| \int_{\mathcal{X}} (u^{k+1} - u^k) \cdot \mathcal{L}_u(u^{k+1} - u^k) \right| \leq C \|\nabla u^{k+1} - \nabla u^k\|_{L^2}^2. \tag{4.18}$$

By applying the the divergence theorem and the inequality (4.18) to (4.11),

$$\mathcal{J}(u^k) - \mathcal{J}(u^{k+1}) \geq \left( \frac{1}{\tau} - C \right) \left\| \nabla u^k - \nabla u^{k+1} \right\|_{L^2}^2. \tag{4.19}$$

Choosing the step size $\tau < \tau_{max} := \frac{1}{2} C^{-1}$ in (4.19), we finish the proof of (4.10).   $\square$

Next, we study the limit of the sequence $\{u^k\}_k \subset V_\lambda$ in proposition 4.3. Since the functional $\mathcal{J}$ is bounded from below, by summing the inequality (4.10), the series

$$\sum_{k=0}^{\infty} \left\| \nabla u^k - \nabla u^{k+1} \right\|_{L^2}^2 < \infty,$$

is convergent and hence $\|\nabla u^k - \nabla u^{k+1}\|_{L^2} \to 0$ as $k \to \infty$. For any $c$-convex function $v$, by the convexity of $\mathcal{J}$ and the divergence theorem,

$$\mathcal{J}(v) - \mathcal{J}(u^k) \geq \int_{\mathcal{X}} \mathcal{J}'(u^k)(v - u^k) = \frac{1}{\tau} \int_{\mathcal{X}} (\Delta u^{k+1} - \Delta u^k)(v - u^k)$$

$$\geq -\frac{1}{\tau} \|\nabla u^k - \nabla u^{k+1}\|_{L^2} \|\nabla v - \nabla u^k\|_{L^2}.$$

By passing the limit $k \to \infty$, we have $\mathcal{J}(v) \geq \lim_{k\to\infty} \mathcal{J}(u^k)$. Since $\{u^k\}_k$ is uniformly bounded in $V_\lambda$, by the Arzela-Ascoli theorem, there exists a subsequence of $\{u^k\}_k$ converges uniformly to some $c$-convex function $\tilde{u}$ under $C^{1,\alpha}$ norm. Hence,

$$\inf_{c\text{-convex}} \mathcal{J}(v) = \lim_{k\to\infty} \mathcal{J}(u^k) = \mathcal{J}(\tilde{u}),$$

As shown in Theorem 4.1, the functional $\mathcal{J}$ has a unique $c$-convex minimizer, which means that $\tilde{u}$ is exactly the unique solution $u$ in Theorem 4.1. Since any convergent subsequence of $\{u^k\}_k$ converges uniformly to $u$, we obtain the following convergence result.

**Theorem 4.4.** *Under the same assumptions of Proposition 4.3, the sequence $\{u^k\}_k$ converges uniformly to $u$ in $\tilde{C}^{1,\alpha}$, $\forall \alpha \in (0,1)$, where $u$ is the solution of the generalized Monge-Ampère equation (3.1).*

It should be noted that Proposition 4.3 utilizes the smoothness of the cost function $c$ on $\mathcal{X} \times \mathcal{Y}$ to get an upper bound for the step size $\tau$. However, Proposition 4.3 still holds if the conditions on $c$ are further relaxed. Notice that all the estimates about $c$ in this section can be restricted to the set

$$\operatorname{graph}(T_u) := \{(x,y) : x \in \mathcal{X}, \, y = T_u(x)\},$$

for $c$-convex function $u \in V_\lambda$. Therefore, the domain of $c$ can be reduced from the original $\mathcal{X} \times \mathcal{Y}$ to

$$G_\lambda := \bigcup_{u \in V_\lambda} \operatorname{graph}(T_u),$$

i.e., the cost function $c$ only needs to be smooth on $G_\lambda$.

For the far-field reflector problem, the cost function $c(x, y) = -\log(1 - x \cdot y)$ on $\mathbb{S}^2$ fails to be smooth on the singularity set

$$\text{sing}(c) := \left\{ (x, y) \in \mathbb{S}^2 \times \mathbb{S}^2 : x = -y \right\}.$$

For this cost function, [26] demonstrates that there exists some constant $\delta > 0$ such that

$$d\left(\text{sing}(c), \text{G}_\lambda\right) > \delta, \tag{4.20}$$

where $d$ is the metric on $\mathbb{S}^2$ and $\delta$ depends only on the parameter $\lambda$. The formula (4.20) means that the reflector cost $c(x, y) = \log(1 - x \cdot y)$ is smooth on $\text{G}_\lambda$. Hence, both Proposition 4.3 and Theorem 4.4 can be extended to the far-field reflector problem.

*Remark* 4.5. The results in Section 3 and Section 4 are discussed for the entire space $\mathbb{S}^n$ or $\mathbb{T}^n$. For the far-field reflector problem, the density functions $f, g$ are supported on $\Omega \subset \mathbb{S}^-$ and $\Omega^* \subset \mathbb{S}^-$, respectively. To reduce the computation, we may solve (3.13) on $\Omega$ instead. An OT boundary condition [3]

$$T_{u^{k+1}}(\partial\Omega) = \partial\Omega^*, \tag{4.21}$$

should be applied to the descent iteration (3.13). The $c$-convexity of $u^{k+1}$ in terms of the space $\Omega \times \Omega^*$ can be ensured by the condition (3.7) together with (4.21), which leads to a diffeomorphism $T_{u^{k+1}}$ from $\Omega$ to $\Omega^*$. It should be noted that the $c$-convexity of $u^{k+1}$ is important for the descent scheme (3.13), since $\mathcal{J}'(u^{k+1})$ can not give a descent direction for $\mathcal{J}$ if $u^{k+1}$ is not $c$-convex.

## 5. Numerical method for the freeform reflector design

In this section, based on the iteration (3.13), we employ the mixed finite element scheme in [29] to solve the far-field reflector problem. The numerical experiments are performed using the open-source finite element framework Firedrake.

We use the notation $\mathcal{M}_h$ to denote the triangular partition of the computational domain, $V_h$ to denote the Lagrange finite element space consisting of continuous piecewise polynomials on $\mathcal{M}_h$, and $\Sigma_h := (V_h)_{3\times 3}$. For simplicity, $\langle \cdot, \cdot \rangle_\Omega$ is used to denote the $L^2$ inner product on $\Omega$.

The far-field design is equivalent to the optimal transport problem with $\mathcal{X} = \mathcal{Y} = \mathbb{S}^2$ and $c(x, y) = -\log(1 - x \cdot y)$. The iterative scheme (3.13) for solving OT is a Poisson equation with respect to $u^{k+1}$, whose weak formulation is written as

$$\left\langle \nabla v, \nabla u^{k+1} \right\rangle = \left\langle \nabla v, \nabla u^k \right\rangle + \tau \left\langle v, r^k \right\rangle, \quad \forall v \in V_h. \tag{5.1}$$

Here, $r^k$ is the residual of the generalized Monge-Ampère equation,

$$r^k(x) := g\left(T_{u^k}(x)\right) \left| \det\left(\nabla T_{u^k}(x)\right) \right| - f^k(x), \tag{5.2}$$

where $T_u$ is the map defined by (3.6) and $f^k(x) := \theta^k f(x)$,

$$\theta^k := \left(\int_{\mathbb{S}^2} f\right)^{-1} \cdot \left(\int_{\mathbb{S}^2} g\left(T_{u^k}\right) \left| \det\left(\nabla T_{u^k}\right) \right| \right).$$

For the far-field reflector problem, $T_u$ is exactly the optical map in (2.3). However, when applying the finite element method, the computational domain $\mathcal{M}_h$ is embedded in $\mathbb{R}^3$. The Jacobian matrix computed here is actually $\nabla_{\mathbb{R}^3} T_u$ instead $\nabla_{\mathbb{S}^2} T_u$ that we need. Hence, we recompute the Jacobian determinant in (5.2) using the basic definition

$$\left| \det\left(\nabla_{\mathbb{S}^2} T_u(x)\right) \right| = \lim_{|U| \to 0} \frac{|T_u(U_x)|}{|U_x|}, \quad U_x \text{ is the neighbourhood of } x. \tag{5.3}$$

The computation of the Jacobian determinant follows the strategies of [29]. For $x \in \mathbb{S}^2$, we parameterize its tangent plane using orthogonal vectors $\{e_1, e_2\}$. The corresponding surface area element $U$ equals to $|\det(e_1, e_2, x)|$. Accordingly, the surface area element of $|T_u(U)|$ is $|\det(v_1, v_2, y)|$, where $v_i = \nabla_{\mathbb{R}^3} T_u(x)\, e_i$, $i = 1, 2$ and $y = T_u(x)$. Constructing the projection matrix,

$$\begin{cases} P(x) := I - xx^{\mathrm{T}}, \\ Q_u(x) := T_u(x)\, x^{\mathrm{T}}, \end{cases} \tag{5.4}$$

we obtain

$$(v_1, v_2, y) = \Big( \nabla_{\mathbb{R}^3} T_u(x) P(x) + Q_u(x) \Big)(e_1, e_2, x).$$

The Jacobian determinant of $T_u : \mathbb{S}^2 \to \mathbb{S}^2$ for $c$-convex $u$ can be given by

$$\Big| \det \big( \nabla_{\mathbb{S}^2} T_u(x) \big) \Big| = \det \Big( \nabla_{\mathbb{R}^3} T_u(x)\, (I - xx^{\mathrm{T}}) + T_u(x)\, x^{\mathrm{T}} \Big), \tag{5.5}$$

which makes it easy to compute $r^k$ in (5.2) under the finite element framework.

In practice, we choose the computational domain as the entire sphere if intensity distributions $f, g$ are supported on complicated subsets of $\mathbb{S}^2$. On the other hand, if the structures of $\Omega := \operatorname{supp} f \subset \mathbb{S}_-^2$ and $\Omega^* := \operatorname{supp} g \subset \mathbb{S}_+^2$ are regular, the computational domain can be chosen as $\Omega$ to improve the efficiency and accuracy.

In this case, the boundary condition (4.21) should be applied to Possion's equation (3.13). Here we use the Neumann boundary condition to realize (4.21) instead, which means that a priori information of $\nabla u^{k+1} \cdot \nu$ is needed, where $\nu$ is the normal of $\partial \Omega$.

Suppose that we already obtain $u^k$ at the $k$-th step. Then $u^k$ defines the corresponding map $T_{u^k}$, which is expected to satisfy the boundary condition $T_{u^k}(\partial \Omega) = \partial \Omega^*$. To achieve this, we project each point $y \in T_{u^k}(\partial \Omega)$ onto the closest point $p$ on $\partial \Omega^*$,

$$p^k(x) := \operatorname{Proj}_{\partial \Omega^*} (T_{u^k}(x)), \quad x \in \partial \Omega, \tag{5.6}$$

where

$$\operatorname{Proj}_{\partial \Omega^*}(y) := \exp_y (H(y) \nabla H(y)), \quad y \in \mathbb{S}^2.$$

Here, the symbol exp denotes the exponential map on the tangent bundle of $\mathbb{S}^2$, and $H(y)$ is the distance function of $\partial \Omega^*$,

$$H(y) := \begin{cases} - \operatorname{dist}(y, \partial \Omega^*), & y \in \Omega^*, \\ + \operatorname{dist}(y, \partial \Omega^*), & y \in \mathbb{S}^2 / \Omega^*. \end{cases} \tag{5.7}$$

We can update the information of $\nabla u^k$ from $p^k(x)$ in (5.6) and denote it as

$$h^{k+1}(x) := \mathcal{T}^{-1} \big( x, p^k(x) \big) = -\nabla_x c \big( x, p^k(x) \big), \quad x \in \partial \Omega. \tag{5.8}$$

For the far-field reflector problem $c(x, y) = -\log(1 - x \cdot y)$, it is equal to

$$h^{k+1}(x) = \frac{p^k(x) - \big( x \cdot p^k(x) \big) x}{1 - x \cdot p^k(x)}, \quad x \in \partial \Omega.$$

The vector function $h^{k+1}$ can be viewed as the information of $\nabla u^{k+1}$ on $\partial \Omega$.

In order to compute the Jacobian determinant precisely, $\nabla_{\mathbb{R}^3} T_{u^k}$ is represented by a tensor-valued variable $\sigma^k \in \Sigma_h$, which solves the linear variational formulation,

$$\big\langle \sigma^k, \chi \big\rangle_\Omega = - \langle \operatorname{div} \chi, T_{u^k} \rangle_\Omega + \langle T_{u^k}, \chi \nu \rangle_{\partial \Omega}, \quad \forall \chi \in \Sigma_h. \tag{5.9}$$

Therefore, the weak form of (3.13) on $\Omega$ can be formulated as

$$\left\langle \nabla u^{k+1}, \nabla v \right\rangle_\Omega = \left\langle \nabla u^k, \nabla v \right\rangle_\Omega + \tau \left\langle r^k, v \right\rangle_\Omega + \left\langle \nu, v(h^k - h^{k+1}) \right\rangle_{\partial\Omega}, \quad v \in V_h. \qquad (5.10)$$

Here, $r^k$ defined by (5.2) is rewritten as

$$r^k(x) = g\big(T_{u^k}(x)\big) \det\left( \sigma^k(x)(I - xx^{\mathrm{T}}) + T_{u^k}(x)\, x^{\mathrm{T}} \right) - f^k(x),$$

$$f^k(x) := \theta^k \cdot f(x),$$

where

$$\theta^k := \left( \int_\Omega f \right)^{-1} \cdot \left( \int_\Omega g(T_{u^k}) \det\left( \sigma^k(I - xx^{\mathrm{T}}) + T_{u^k}\, x^{\mathrm{T}} \right) \right).$$

The iteration should be terminated when the functional $\mathcal{J}(u^k)$ stops decreasing. However, it is complicated to compute $\mathcal{J}(u^k)$ due to the computation of $u^c(y)$. Instead, we use the $L^2$ norm of $r^k$,

$$\|r^k\|_2 := \left( \int_\Omega |r^k(x)|^2 \,\mathrm{d}x \right)^2,$$

as the stop criterion. The numerical method for solving the far-field reflector problem on $\Omega$ is summarized in Algorithm 1.

---

**Algorithm 1** FEM for the far-field reflector design on $\Omega$.

---

Given the initial value $u^0$ and step size $\tau$.
**while** $\|r^{k-1}\|_2 < \|r^{k-2}\|_2$ **do**
    Solving the linear variational problem (5.9) to obtain $\sigma^k$.
    Evaluating the composite function $g(T_{u^k})$.
    Solving the weak formulation of Poisson's equation (5.10) to obtain $u^{k+1}$.
    Computing the error $\|r^k\|_2$.
    $k := k + 1$.
**end while**

---

## 6. Numerical experiments

In this section, we present the numerical experiments for the far-field reflector problem using the algorithm in Section 5. We define the subdomain $\mathbb{S}^2_{-\theta} \subset \mathbb{S}^-$ as

$$\mathbb{S}^2_{-\theta} := \left\{ x \in \mathbb{S}^2 : x \cdot e_{\boldsymbol{z}} \leq -\cos(\theta),\ \text{where } e_{\boldsymbol{z}} = (0,0,1)^{\mathrm{T}} \right\}, \quad \theta \in \left( 0, \frac{\pi}{2} \right),$$

and $\mathbb{S}^2_{+\theta} \subset \mathbb{S}^+$ as

$$\mathbb{S}^2_{+\theta} := \left\{ x \in \mathbb{S}^2 : x \cdot e_{\boldsymbol{z}} \geq \cos(\theta),\ \text{where } e_{\boldsymbol{z}} = (0,0,1)^{\mathrm{T}} \right\}, \quad \theta \in \left( 0, \frac{\pi}{2} \right).$$

We present several numerical experiments by designing different target intensities $g$. The source intensity

$$f(x) = \begin{cases} 1, & x \in \mathbb{S}^2_{-\pi/4}, \\ 0, & \text{else}, \end{cases} \qquad (6.1)$$

supported on $\Omega = \mathbb{S}^2_{-\pi/4}$ is fixed. In fact, the computational stability of the numerical method depends more on the target $g$ than the source $f$. In experiments, $\mathcal{M}_h$ is the second-order mesh discretized on $\Omega = \mathbb{S}^2_{-\pi/4}$. The finite element space $V_h$ is of second order.

Linear equations in Algorithm 1 are solved by incomplete LU factorization preconditioned CG method, where the constraint $\int_\Omega u^{k+1} = 0$ is handled by the constant null space of the Krylov solver.

In this section, $\|r\|_2 := \|r(x)\|_{L^2}$ is used to denote the $L^2$ norm of the residual. The illumination is simulated using the ray-tracing method. We plot the ray-traced image $T_{u_h \#} f$ using the kernel density function in Matlab, where $u_h$ is the numerical solution of the problem, and the notation $\#$ is defined in (2.1). To ensure the convergence of the method, the initial value $u^0$ should be a $c$-convex function, which is generally zero. If the target $g$ is far away from $T_{u^0 \#} f$, we may choose a middle value $g_{mid}$ and solve the reflector problem with the target $g_{mid}$. Then, the solution of this problem can be used as the initial value for the reflector problem with the target $g$.

**Example 6.1. Smooth target.** In this example, we design a freeform reflector to produce the smooth target intensity $g$ in the far-field. Here, the intensity $g$ is compactly support on $\mathbb{S}^2_{+\pi/4}$, shown Figure 3(a). The source intensity $f$ is given by (6.1). The initial function $u^0$ in Algorithm 1 is set to be zero. We choose the step size $\tau = 0.5$ and the finite element mesh size $h = 9.82 \times 10^{-3}$.



FIGURE 2. Results of Example 6.1. (a) Numerical solution $u_h$ on $\mathbb{S}^2_{-\pi/4}$, where the mesh size $h = 9.82 \times 10^{-3}$. (b) Radial distance function $\rho_h := e^{-u_h}$. (c) Magnitude $|\nabla \rho_h|$. (d) Vector field $\nabla \rho_h$, where the colors correspond to the values of $|\nabla \rho_h|$.

Figure 2(a) and (b) demonstrate the numerical solution $u_h$ and its gradient $\nabla u_h$, respectively. In Figure 3(b), we presents the ray-traced result $g_h := T_{u_h \#} f$ using the computed $u_h$ in Figure 2(a). It is difficult to distinguish the difference between $g$ and $g_h$ from Figure 3(a) and (b). Further, we plot the absolute error $|g - g_h|$ in Figure 3(c) to evaluate the
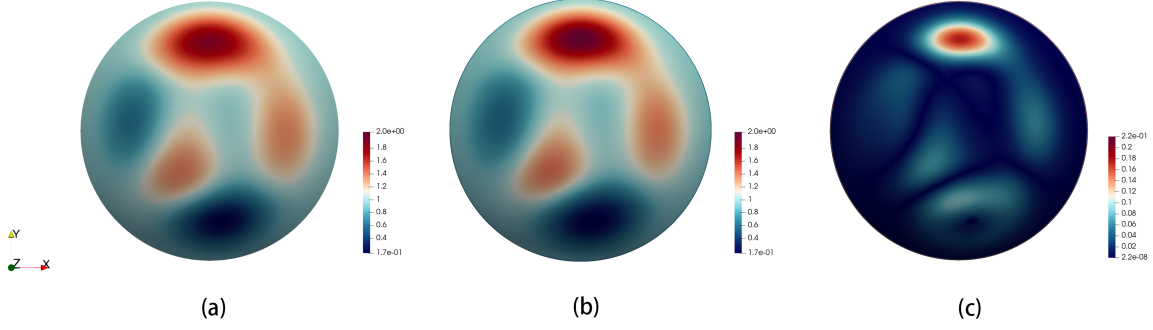
FIGURE 3. Ray-traced image of Example 6.1. (a) Far-field target intensity $g$ on $\mathbb{S}^2_{+\pi/4}$. (b) Far-field ray-traced intensity $g_h$ using the numerical solution $u_h$, where $h = 9.82 \times 10^{-3}$. (c) Absolute error $|g - g_h|$.

difference between $g$ and $g_h$. The error is relatively small in most areas, and the maximal absolute error 0.2 occurs near the peak of the intensity $g$. Although the intensity $g$ is smooth, the high contrast of $\frac{\max g}{\min g} \approx 12$ makes the condition number of $\nabla T_{u_h}$ small, decreasing the accuracy of the numerical solution.



FIGURE 4. Convergence curves of Example 6.1 at different discretization levels, residual $\|r\|_2$ versus iterations.

Figure 4 shows the changes of residual value $\|r\|_2$ over iterations for mesh levels $N = 20, 40, 80, 160$, which correspond to the mesh sizes $h = 3.93 \times 10^{-2}, 1.96 \times 10^{-2}, 9.82 \times 10^{-3}, 4.9 \times 10^{-3}$. Here $N$ refers to the number of mesh cells along the geodesic radius of $\mathbb{S}^2_{-\pi/4}$. The algorithms for $N = 20, 40, 80, 160$ stop after 11, 21, 45, and 47 iterations, respectively.

**Example 6.2. Off-axis target**. In this example, the far-field target intensity $g$ is

$$
g(x) = \begin{cases} 1, & \text{if } x \cdot q \geq \cos\left(\frac{\pi}{4}\right), \text{ where } q = \left(0, -\sin\left(\frac{\pi}{8}\right), \cos\left(\frac{\pi}{8}\right)\right)^{\mathrm{T}}, \\ 0, & \text{else}, \end{cases} \tag{6.2}
$$

which corresponds to the blue region in Figure 5(a). The source target $f$ defaults to (6.1), corresponding to the orange region in Figure 5(a). The step size $\tau = 0.5$ and the mesh size $h = 9.82 \times 10^{-3}$.
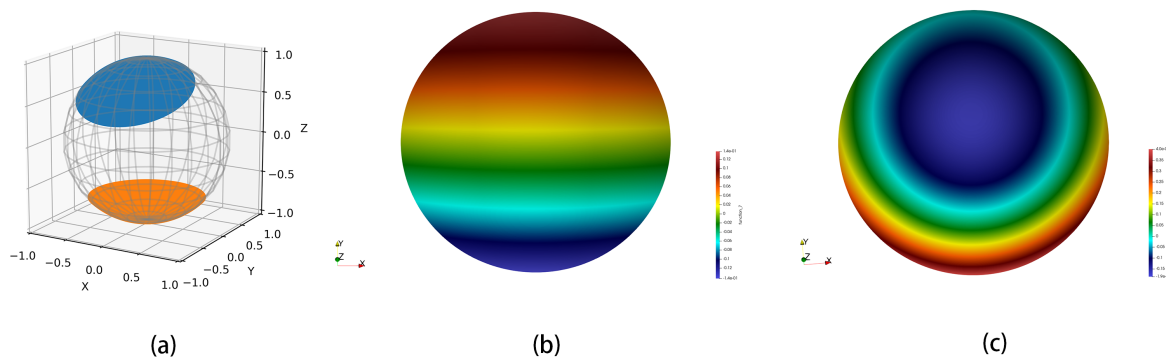
FIGURE 5. Results of Example 6.2. (a) The source intensity $f$ (orange) and the far-field target intensity $g$ (blue) on $\mathbb{S}^2$. (b) Numerical solution $u_h$ on $\mathbb{S}^2_{-\pi/4}$, where the mesh size $h = 9.82 \times 10^{-3}$. (c) Numerical solution $\varphi_h$ with $h = 9.82 \times 10^{-3}$.
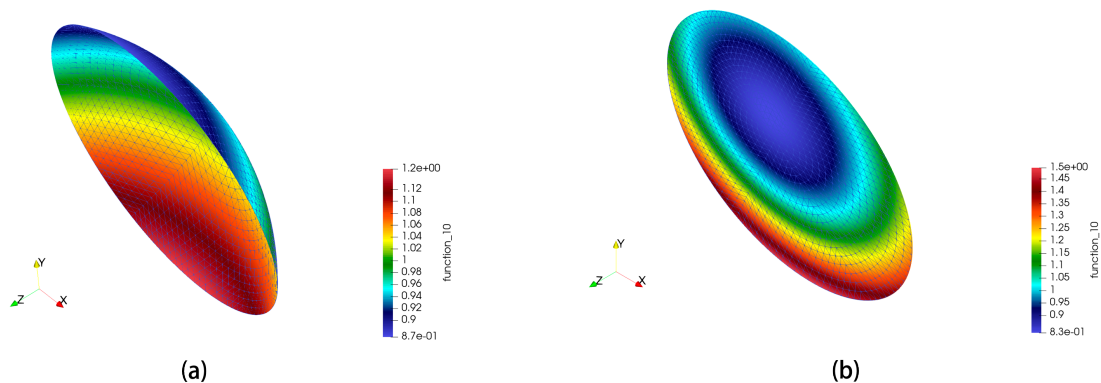


FIGURE 6. Results of Example 6.2. (a) The reflecting surface $\Gamma = \{x\rho_h(x), \ x \in \mathbb{S}^2_{-\pi/4}\}$, where $\rho_h = e^{-u_h}$. (b) The reflecting surface $\Gamma = \{x\rho_h(x), \ x \in \mathbb{S}^2_{-\pi/4}\}$, where $\rho_h = e^{\varphi_h}$. The colors of the surfaces correspond to the values of $\rho_h$.

In this example, we solve the far-field problem using two different formulations mentioned in Theorem 2.1 and Theorem 2.2, i.e., the optimal transport problem with the cost function $c(x,y) = -\log(1 - x \cdot y)$ and $c(x,y) = \log(1 - x \cdot y)$. Here, their numerical solutions are denoted by $u_h$ and $\varphi_h$, which are shown in Figure 5 (b) and (c), respectively. The corresponding reflecting surfaces are shown in Figure 6.

In experiments, we initialize

$$u^0 = 0,$$

$$\varphi^0 = \log\left(\frac{b}{x \cdot e_z}\right), \ e_z = (0, 0, 1)^{\mathrm{T}},$$

and $b$ is some negative constant such that $\int_{\mathbb{S}^2_{-\pi/4}} \varphi^0 = 0$. In particular, $T_{u^0 \#}f = T_{-\varphi^0 \#}f = f(-x)$, i.e., the uniform distribution on $\mathbb{S}^2_{+\pi/4}$. The center of $f(-x)$ forms an angle of $\frac{\pi}{8}$ with the center of $g$, increasing the difficulty of the boundary match.
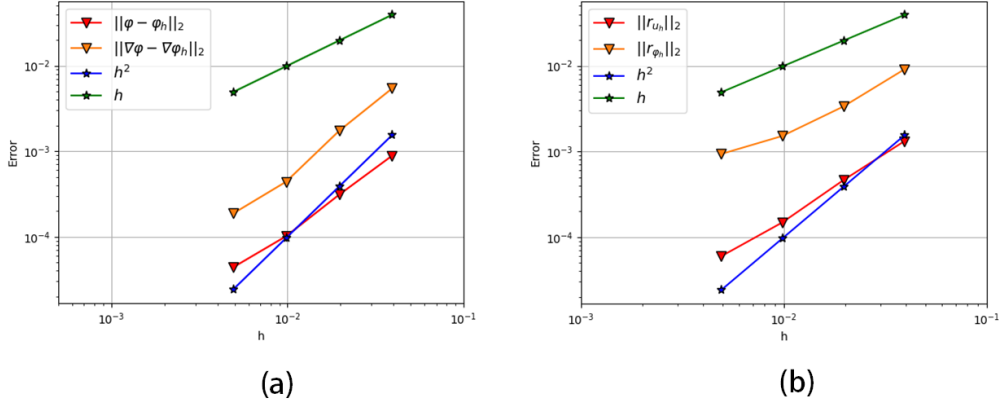
FIGURE 7. Error plots for the numerical solutions $u_h, \varphi_h$ of Example 6.2.
(a) $\|\varphi - \varphi_h\|_2$ and $\|\nabla\varphi - \nabla\varphi_h\|_2$ versus the mesh size $h$, where $\varphi$ is the exact
solution. (b) Residual values $\|r_{u_h}\|_2$ and $\|r_{\varphi_h}\|_2$ versus $h$, where $r_{u_h}$ and $r_{\varphi_h}$
are the residuals of $u_h$ and $\varphi_h$, respectively.

We remarked that for this example, OT with the cost $c(x, y) = \log(1 - x \cdot y)$ has the
exact solution $\varphi$

$$\varphi(x) = \log\left(\frac{b}{x \cdot p}\right), \qquad p = \left(\frac{q_1}{\sqrt{2(q_3 + 1)}}, \frac{q_2}{\sqrt{2(q_3 + 1)}}, \sqrt{\frac{q_3 + 1}{2}}\right)^{\mathrm{T}}, \qquad x \in \mathbb{S}^2_{-\pi/4},$$

where the vector $q = (q_1, q_2, q_3)^{\mathrm{T}}$ is given in (6.2) and $b$ is an arbitrary negative constant,
which is often taken to make $\int_{\mathbb{S}^2_{-\pi/4}} \varphi = 0$. The exact reflector surface of this case is a plane
in $\mathbb{R}^3$, as shown in Figure 6(b). We plot the errors $\|\varphi - \varphi_h\|_2$ and $\|\nabla\varphi - \nabla\varphi_h\|_2$ versus
the mesh size $h$ in Figure 7(a). Meanwhile, the changes of the residual values $\|r_{u_h}\|_2$ and
$\|r_{\varphi_h}\|_2$ with respect to the mesh size $h$ have been shown in Figure 7(b). It is obvious that
the solution $u_h$ has higher accuracy than $\varphi_h$.

**Example 6.3. Discontinuous target with singular background**. Different from the
previous examples, we set the far-field illumination area to be a far-field plane $P_\infty$ instead
of a far-field sphere. In addition, the target intensity $g$ is supported on the square area of
$P_\infty$, which has an image of the letter "A" as shown in Figure 10(a). For the visualization
of the image on the plane, $P_{0.5} := \{x = (x_1, x_2, x_3)^{\mathrm{T}} \in \mathbb{R}^3 : x_3 = 0.5\}$ is used to substitute
$P_\infty$. We use $T_u^{\mathrm{P}}$ to denote the reflecting map from the source $\mathbb{S}^2$ to the target plane $P_{0.5}$.

Utilizing the stereographic projection, we can project $g$ on $P_{0.5}$ onto a new function $\hat{g}$
on $\mathbb{S}^2$, where $\hat{g}$ defined by $\hat{g} := \mathcal{S}g$,

$$\mathcal{S} : L^1(P_{0.5}) \to L^1(\mathbb{S}^2_{-\pi/4}),$$
$$g \to \hat{g}, \qquad\qquad \text{where}$$

$$\hat{g}(x) = \frac{(X^2 + Y^2 + 0.5^2)^{1.5}}{0.5} g(X, Y), \quad x = (x_1, x_2, x_3)^{\mathrm{T}} \in \mathbb{S}^2_{-\pi/4},$$

$$\text{where} \quad \begin{cases} X = 0.5\dfrac{\sqrt{x_1^2 + x_2^2}}{x_3}\cos\left(\tan^{-1}\left(\dfrac{x_2}{x_1}\right)\right), \\ Y = 0.5\dfrac{\sqrt{x_1^2 + x_2^2}}{x_3}\sin\left(\tan^{-1}\left(\dfrac{x_2}{x_1}\right)\right). \end{cases}$$

Figure 8(a) shows the image of $\hat{g}$. We can solve the far-field reflector problem for the source $f$ in (6.1) and the target $\hat{g}$. The obtained solution $u$ can produce the illumination pattern of Figure 10(a) on the far-field plane, and the illumination $T^P_{u\#}f$ is related to $T_{u\#}f$ by $T^P_{u\#}f = \mathcal{S}^{-1}(T_{u\#}f)$.

Both the shape of $\operatorname{supp} g$ and the discontinuity of $g$ increase the difficulty of solving for $u$ efficiently. We set the step size $\tau = 0.3$ and the initial value $u^0 = 0$. Due to the discontinuity of the target $g$, the algorithm stops after 12 iterations for $h = 9.82 \times 10^{-3}$, with the residual value $r_{u_h}$ equal to $4 \times 10^{-2}$.
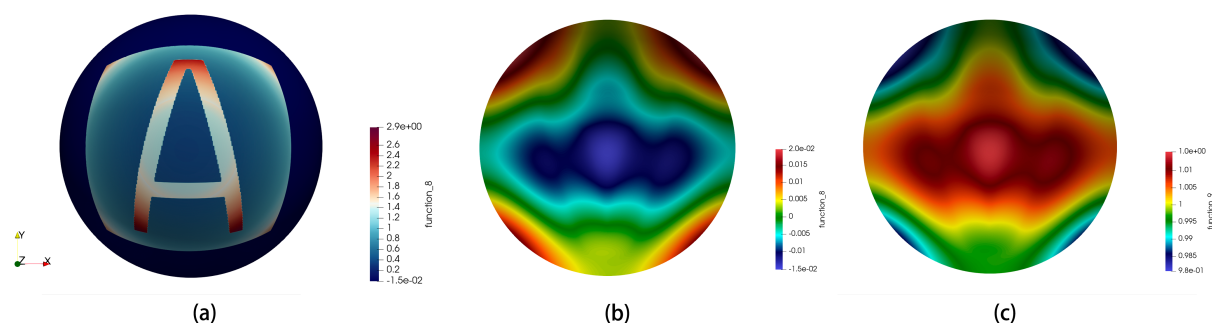


FIGURE 8. (a) Far-field target intensity function $\hat{g}$ on $\mathbb{S}^2_+$, (b) numerical solution $u_h$ on $\mathbb{S}^2_{-\pi/4}$ for $h = 9.82 \times 10^{-3}$ and (c) radial distance function $\rho_h := e^{-u_h}$.
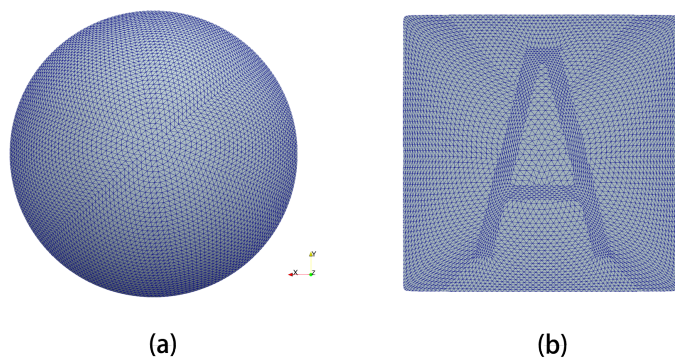


FIGURE 9. Results of Example 6.3. (a) Finite element mesh $\mathcal{M}_h$ with $h = 1.96 \times 10^{-2}$. (b) New mesh obtained by the far-field mapping $T^P_{u_h}\mathcal{M}_h$, $h = 1.96 \times 10^{-2}$.

Figure 8(b) and (c) show the numerical solution $u_h$ of the far-field problem and its radial distance function $\rho_h$, where $h = 9.82 \times 10^{-3}$. Figure 9(a) and (b) respectively show the finite element mesh $\mathcal{M}_h$ discretizing $\mathbb{S}^2_{-\pi/4}$ and the mesh obtianed by $T^P_{u_h}\mathcal{M}_h$, where $h = 1.96 \times 10^{-2}$.

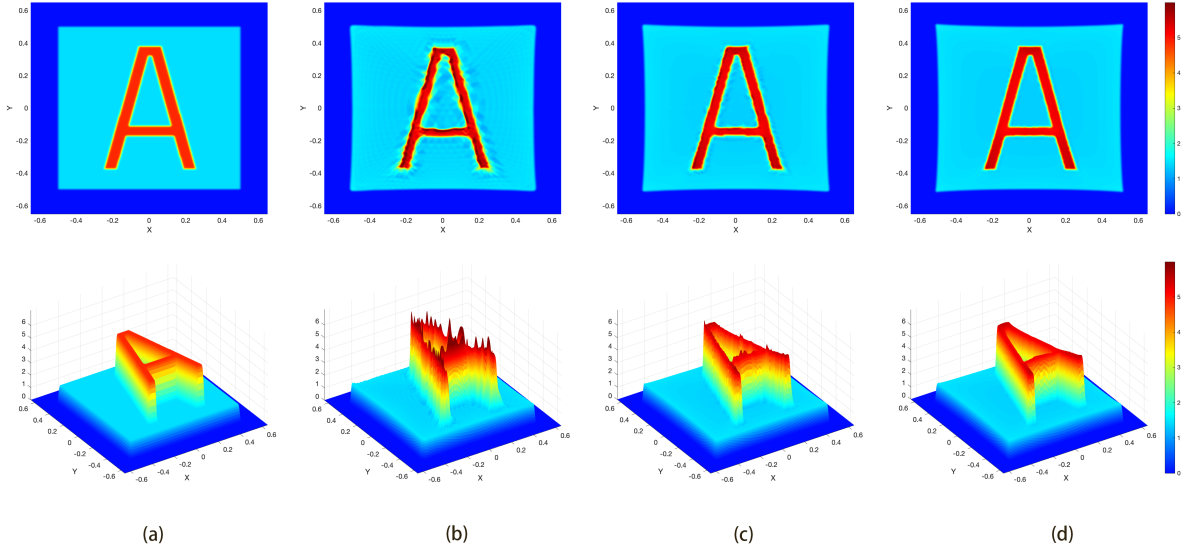FIGURE 10. Ray-traced images $g_h := T^{\mathrm{P}}_{u_h} \# f$ on the plane, using the results $u_h$ of Example 6.3. From left to right: (a) target intensity $g$ on $P_{0.5}$, (b) ray-traced intensity $g_h$ with $h = 3.93 \times 10^{-2}$, (c) ray-traced intensity $g_h$ with $h = 1.96 \times 10^{-2}$, (d) ray-traced intensity $g_h$ with $h = 9.82 \times 10^{-3}$. The first row is the image on the 2D plane, and the second row is its 3D view.

In Figure 10, the far-field ray-traced images $g_h := T^{\mathrm{P}}_{u_h} \# f$ on the target plane $P_{0.5}$ are shown for the mesh sizes (b) $h = 3.93 \times 10^{-2}$, (c) $h = 1.96 \times 10^{-2}$, (d) $h = 9.82 \times 10^{-3}$. The ray-traced illumination on the target plane is increasingly clear as the mesh is refined.

## 7. Conclusions

This work introduces a fast algorithm derived from the optimal transportation theory to address the far-field reflector problem. The method employs the Sobolev gradient descent on the reduced objective functional of OT, achieving a balance between computational stability and efficiency. Theoretical analysis verifies the superior convergence properties of the descent scheme, while experimental results demonstrate that the proposed method performs exceptionally well even in singular cases.

Future work may focus on extending this research to the challenging near-field reflector problem. Another promising direction is to apply the $c$-transform technique to enhance the stability of the method, enabling it to handle highly contrasting singular cases [14] more effectively.

## References

[1] G. Awanou, Standard finite elements for the numerical resolution of the elliptic Monge–Ampére: classical solutions, IMA Journal of Numerical Analysis 35.3 (2015): 1150-1166.

[2] G. Bao, and Y. Zhang, An optimal transport approach for 3D electrical impedance tomography, Inverse Problems 40.12 (2024): 125006.

[3] J. D. Benamou, B. D. Froese, A. M. Oberman, Numerical solution of the optimal transportation problem using the Monge–Ampère equation, Journal of Computational Physics 260 (2014): 107-126.

[4] J. D. Benamou, W. L. Ijzerman and G. Rukhaia, An entropic optimal transport numerical approach to the reflector problem, Methods and Applications of Analysis (2020).

[5] R. J. Berman, The Sinkhorn algorithm, parabolic optimal transport and geometric Monge–Ampére equations, Numerische Mathematik 145.4 (2020): 771-836.

[6] C. Bösel, and H. Gross, Single freeform surface design for prescribed input wavefront and target irradiance, J. Opt. Soc. Am. A 34.9 (2017): 1490-1499.

[7] Y. Brenier, Polar factorization and monotone rearrangement of vector-valued functions, Comm. Pure Appl. Math. 44 (1991), 375–417.

[8] K. Brix, Y. Hafizogullari, and A. Platen, Solving the Monge–Ampère equations for the inverse reflector problem, Mathematical Models and Methods in Applied Sciences 25.05 (2015): 803-837.

[9] D. A. Bykov, L. L. Doskolovich, A. A. Mingazov, et al., Linear assignment problem in the design of freeform refractive optical elements generating prescribed irradiance distributions, Optics Express 26.21 (2018): 27812-27825.

[10] L. A. Caffarelli, Interior $W^{2,p}$ estimates for solutions of the Monge-Ampere equation, Annals of Mathematics 131.1 (1990): 135-150.

[11] L. A. Caffarelli, C. E. Gutiérrez, and Q. Huang, On the regularity of reflector antennas, Annals of Mathematics (2008): 299-323.

[12] P. M. M. De Castro, Q. Mérigot, and B. Thibert, Far-field reflector problem and intersection of paraboloids, Numerische Mathematik 134 (2016): 389-411.

[13] K. Desnijder, P. Hanselaer, and Y. Meuret, Ray mapping method for off-axis and non-paraxial freeform illumination lens design, Optics letters 44.4 (2019): 771-774.

[14] L. L. Doskolovich, E. V. Byzov, A. A. Mingazov, et al., Supporting quadric method for designing freeform mirrors that generate prescribed near-field irradiance distributions, Photonics. Vol. 9. No. 2. MDPI, 2022.

[15] L. C. Evans, Partial Differential Equations, Vol. 19. American Mathematical Society, 2022.

[16] Z. Feng, B. D. Froese, and R. Liang, Freeform illumination optics construction following an optimal transport map, Applied Optics 55.16 (2016): 4301-4306.

[17] X. Feng, and M. Neilan, Mixed finite element methods for the fully nonlinear Monge–Ampère equation based on the vanishing moment method, SIAM Journal on Numerical Analysis 47.2 (2009): 1226-1250.

[18] F. R. Fournier, W. J. Cassarly, and J. P. Rolland, Fast freeform reflector generation using source-target maps, Optics Express 18.5 (2010): 5295-5304.

[19] T. Glimm, and V. Oliker, Optical design of single reflector systems and the Monge–Kantorovich mass transfer problem, Journal of Mathematical Sciences 117.3 (2003): 4096-4108.

[20] B. F. Hamfeldt, and A. G. Turnquist, Convergent numerical method for the reflector antenna problem via optimal transport on the sphere, J. Opt. Soc. Am. A 38.11 (2021): 1704-1713.

[21] Y. H. Kim, J. Streets, and M. Warren, Parabolic optimal transport equations on manifolds, International Mathematics Research Notices 2012.19 (2012): 4325-4350.

[22] S. A. Kochengin, and V. I. Oliker, Determination of reflector surfaces from near-field scattering data II. Numerical solution, Numerische Mathematik 79.4 (1998): 553-568.

[23] S. A. Kochengin, and V. I. Oliker, Computational algorithms for constructing reflectors, Computing and Visualization in Science 6.1 (2003): 15-21.

[24] G. Loeper, and F. Rapetti, Numerical solution of the Monge–Ampère equation by a Newton's algorithm, Comptes Rendus Mathematique 340.4 (2005): 319-324.

[25] G. Loeper, On the regularity of solutions of optimal transportation problems, Acta Mathematica 202.2 (2009): 241-283.

[26] G. Loeper, Regularity of optimal maps on the sphere: The quadratic cost and the reflector antenna, Archive for Rational Mechanics and Analysis 199 (2011): 269-289.

[27] X. N. Ma, N. S. Trudinger, and X. J. Wang, Regularity of potential functions of the optimal transportation problem, Archive for Rational Mechanics and Analysis 177 (2005): 151-183.

[28] R. J. McCann, Polar factorization of maps on Riemannian manifolds, Geometric & Functional Analysis GAFA 11.3 (2001): 589-608.

[29] A. T. McRae, C. J. Cotter, and C. J. Budd, Optimal-transport-based mesh adaptivity on the plane and sphere using finite elements, SIAM Journal on Scientific Computing 40.2 (2018): A1121-A1148.

[30] J. Meyron, Q. Mérigot, and B. Thibert, Light in power: a general and parameter-free algorithm for caustic design, ACM Transactions on Graphics (TOG) 37.6 (2018): 1-13.

[31] J. Neuberger, Sobolev gradients and differential equations, Springer Science & Business Media, 2009.

[32] C. R. Prins, R. Beltman, J. H. M. ten Thije Boonkkamp, W. L. IJzerman, and T. W. Tukker, A least-squares method for optimal transport using the Monge-Ampère equation, SIAM Journal on Scientific Computing 37.6 (2015): B937-B961.

[33] F. Rathgeber, D. A. Ham, L. Mitchell, et al., Firedrake: Automating the finite element method by composing abstractions, ACM Trans. Math. Software, 43 (2017), 24.

[34] L. B. Romijn, J. H. M. ten Thije Boonkkamp, M. J. H. Anthonissen, and W. L. IJzerman, An iterative least-squares method for generated Jacobian equations in freeform optical design, SIAM Journal on Scientific Computing 43.2 (2021): B298-B322.

[35] M. M. Sulman, J. F. Williams, and R. D. Russell, An efficient approach for the numerical solution of the Monge–Ampère equation, Applied Numerical Mathematics 61.3 (2011): 298-307.

[36] C. Villani, Optimal transport: old and new. Vol. 338. Berlin: Springer, 2009.

[37] X. J. Wang, On the design of a reflector antenna, Inverse Problems 12.3 (1996): 351.

[38] X. J. Wang, On the design of a reflector antenna II, Calculus of Variations and Partial Differential Equations 20.3 (2004): 329-341.

[39] R. Wu, L. Xu, P. Liu, Y. Zhang, Z. Zheng, H. Li, and X. Liu, Freeform illumination design: a nonlinear boundary problem for the elliptic Monge–Ampére equation, Optics Letters 38.2 (2013): 229-231.

SCHOOL OF MATHEMATICAL SCIENCES, ZHEJIANG UNIVERSITY, HANGZHOU 310027, CHINA.
*Email address*: baog@zju.edu.cn

SCHOOL OF MATHEMATICAL SCIENCE, PEKING UNIVERSITY, BEIJING 100871, CHINA.
*Email address*: yixuan@pku.edu.cn