# CAUSAL TRANSPORT IN DISCRETE TIME AND APPLICATIONS

JULIO BACKHOFF, MATHIAS BEIGLBÖCK, YIQING LIN, AND ANASTASIIA ZALASHKO

ABSTRACT. Loosely speaking, causal transport plans are a relaxation of adapted processes in the same sense as Kantorovich transport plans extend Monge-type transport maps. The corresponding causal version of the transport problem has recently been introduced by Lassalle. Working in a discrete time setup, we establish a dynamic programming principle that links the causal transport problem to the transport problem for general costs recently considered by Gozlan et al. Based on this recursive principle, we give conditions under which the celebrated Knothe-Rosenblatt rearrangement can be viewed as a causal analogue to the Brenier's map. Moreover, these considerations provide transport-information inequalities for the nested distance between stochastic processes pioneered by Pflug and Pichler, and so serve to gauge the discrepancy between stochastic programs driven by different noise distributions.

**Keywords:** Optimal transport, causality, nested distance, general transport costs, Knothe-Rosenblatt rearrangement, transport inequalities.

**AMS subject classifications.** 90C15,60G70,39B62

## 1. INTRODUCTION

In this article we consider the optimal transport problem between two discrete-time stochastic processes under the so-called causality constraint, highlighted recently by the work of Lassalle in [Las15] in a more general setting. A transport plan between two processes is said to be *causal* if, from an observed trajectory of the first process, the "mass" can be split at each moment of time into the second process only based on the information available up to that time. It is illustrative to think of the deterministic case (i.e. when there is no splitting of mass); such a causal plan is then an actual mapping which is further adapted, and so the relationship between causal plans and adapted processes is the same as between classical transport plans (Kantorovich) and transport maps (Monge).

The idea of imposing a "causality" constraint on a transport plan between laws of processes seems to go back to the Yamada-Watanabe criterion for stochastic differential equations [YW71]. Under the name "compatibility" the same type of constraint was introduced in Kurz [Kur07]. In this article we will also link these objects to the notion of *nested distance*, whose systematic investigation was initiated by Pflug [Pfl09] and Pflug–Pichler [PP12, PP14, PP15], and had a precursor in the "Markov-constructions" studied by Rüschendorf [Rüs85]. Roughly, the nested distance is defined through a problem of optimal transport over plans which are *bicausal*, this notion being the symmetrized analogue of causality. Interestingly, [Rüs85] and [PP14] established a recursive formulation for the problem, and [PP12, PP14] further obtained a dual formulation for the nested distance. Moreover, Pflug–Pichler [PP12] applied these considerations to the practical problem of reducing the complexity of multistage stochastic programs, by showing that the difference between the optimal value of a program w.r.t. two different noise distributions is dominated by the nested distance between them. We refer to the books [RS03, SDR14, PP14] for a detailed account on stochastic programming.

A systematic treatment and use of causality as an interesting property of abstract transport plans and their associated optimal transport problems was first made by Lassalle in [Las13] in the general context of Polish spaces (then updated in [Las15]). As an application the author considers the Wiener space setting of the problem and establishes that weak solutions to Brownian-motion-driven stochastic differential equation can be conceived as causal transport plans between the Wiener measure and a target measure, and finds that such plans are automatically bicausal and optimal for a Cameron-Martin-type cost. He then explores functional inequalities in Wiener space (first obtained by [FÜ04]) by means of this

method. We stress that the main motivation for our article comes from this connection with stochastic analysis, and our goal is to deepen the understanding of the causal transport problem by looking at the discrete-time setup; the continuous-time counterpart/extension of our results is a work in progress. This motivation implies that for some results we are content with assuming independence of marginals or of increments for some of the processes we look at. It also means that we often seek to show the robustness of some of the particular phenomena obtained by Lassalle: for instance by studying when the causal and bicausal problems coincide/differ, or by showing that functional inequalities are also prevalent in our setting.

The core subject of our article is the causal optimal transport problem, which consists in finding the cheapest causal transport plan from a given source measure (process) to a target one, with respect to a certain cost function on the product space. Since causality can be easily characterized as a linear constraint on transports, we can embed this problem in the class of optimal transport problems under (infinitely many) linear constraints, as considered by Zaev [Zae15] and [BG14]; this line of reasoning has already been applied in the literature, for example in the development of martingale optimal transport (see [HN12, BHLP13, GHLT14, DS14]). In this way, we obtain conditions for the existence of an optimal causal transport and identify a dual formulation for the problem, further establishing no-duality-gap; this is the content of our Theorem 2.5. By studying the conditional distributions of causal transports, we are able to tackle many instances of the causal optimal transport problem by means of a recursion, which we call the dynamic programming principle (DPP in short); see Theorem 2.6. The appeal of these recursions is that instead of one "multi-dimensional" transport problem over causal plans, we obtain recursively several "one-dimensional" problems, each one of them a "general" (i.e. non-linear) transport problem as introduced recently in [GRST15]. In this way, we reduce the dimensionality of the problem at the expense of introducing non-linearities.

In Theorem 2.7 we establish that the Knothe-Rosenblatt rearrangement [Kno57] (also known as multi-dimensional quantile transform / increasing triangular transformation) is causal optimal for the squared euclidean distance (or more generally convex and separable in the sense of (2.10) below) if the source measure is of product type. This setting is relevant e.g. in the theory of functional equalities, and such result can be extended to the case when the source measure has independent increments and the cost is suitably modified (see Corollary 2.8) which is a set-up relevant to stochastic analysis. The key here is to first identify the mentioned rearrangement as a bicausal optimizer, which we do in Proposition 5.3 generalizing a corresponding result in [Rüs85], and then to prove that under the given assumptions the values of the causal and bicausal problems coincide. If further the source measure has absolutely continuous marginals, then the Knothe-Rosenblatt rearrangement is of Monge-type. Hence the Knothe-Rosenblatt rearrangement can be viewed as a causal version of the classical Brenier map. In our opinion this result adds to the appeal of this rearrangement, which is in any way widely used in analysis, statistics, and operations research in the context of scenario generation.

We also further the understanding of the relation between nested distance and multistage stochastic programming established in [Pfl09, PP12, PP14]. First we show that many stochastic programs are concave/convex along what we define as lexicographic-displacement interpolations, in analogy to the concepts of displacement interpolation and displacement convex functionals in classical optimal transport theory (see [Vil03, chapter 5] or McCann [McC97]), where the role of Brenier's map is taken by the Knothe-Rosenblatt rearrangement. Moreover, we give conditions under which the nested distance of "order one", which gauges the discrepancy between stochastic programs driven by different noise distributions and a common Lipschitz-cost criterion, can itself be assessed by the square root of the relative entropy between such processes. In other words, we establish a transport-information inequality for this nested distance of "order one". This means that the discrepancy between such stochastic programs can be simply gauged by an entropy, which is easier to compute in practice than the nested distance itself. We shall also have occasion to further highlight the connection between causality and functional inequalities when, in Section 5.1, we establish Talagrand's celebrated $\mathcal{T}_2$ inequality (see [Tal96]) for the standard Gaussian measure by interpreting the author's "tensorization/inductive trick" as an instance of our recursions; we refer the reader to [Led01, GL10] for an account on functional/geometric inequalities and the related concept of concentration of measure.

The article is organized as follows. In Section 2 we introduce the setting and collect our main results; Theorem 2.5 on the attainability and duality for the causal transport problem (established in Section 3),

Theorem 2.6 on the recursive formulation of the problem (what we call the DPP, whose proof is given in Section 4), Theorem 2.7 on the identification of the Knothe-Rosenblatt rearrangement as a causal optimizer (established in Section 5), and finally Theorem 2.9 on the bicausal transport information inequality which is further explored in Section 6 along with other connections to stochastic programming. In Section 7 we present some counterexamples cited throughout the paper.

**Notation:** For a product of sets $\mathcal{X} \times \mathcal{Y}$ we denote by $p^1, p^2$ the projection onto the first resp. second coordinate. The pushforward of a measure $\gamma$ by a map $M$ is denoted $M_*\gamma$. We denote by $\gamma^x, \gamma^y$ the regular kernels of a measure $\gamma$ on $\mathcal{X} \times \mathcal{Y}$ w.r.t. its first and second coordinate respectively. Analogous notation extends to products of more than two spaces. On $\mathbb{R}^N \times \mathbb{R}^N$ we denote by $(x_1, \ldots, x_N)$ the first half and $(y_1, \ldots, y_N)$ the second half of the coordinates, and we convene that for $\gamma$ a probability in $\mathbb{R}^N \times \mathbb{R}^N$ (respect. $\eta$ on $\mathbb{R}^N$), $\gamma^{x_1, \ldots, x_t, y_1, \ldots, y_t}$ (respect. $\eta^{x_1, \ldots, x_t}$) denotes the two-dimensional measure on $(x_{t+1}, y_{t+1})$ (respect. one-dimensional measure on $x_{t+1}$) given by regular disintegration of $\gamma$ w.r.t. $(x_1, \ldots, x_t, y_1, \ldots, y_t)$ (respect. $\eta$ w.r.t. $(x_1, \ldots, x_t)$). Also, a statement like "for $\gamma$-a.e. $x_1, \ldots, x_t, y_1, \ldots, y_t$" or "for $\eta$-a.e. $x_1, \ldots, x_t$" is meant to denote respectively "almost-everywhere" with respect to the projections of $\gamma$ onto $x_1, \ldots, x_t, y_1, \ldots, y_t$ or $\eta$ onto $x_1, \ldots, x_t$.

## 2. MAIN RESULTS

Let $\mathcal{X}$ and $\mathcal{Y}$ be closed subsets of $\mathbb{R}^N$ and take $\mathcal{F}^{\mathcal{X}}$ and $\mathcal{F}^{\mathcal{Y}}$ the filtrations generated by the coordinate processes (i.e. $\mathcal{F}_t^X$ is the smallest $\sigma$-algebra s.t. $x \in \mathcal{X} \mapsto (x_1, \ldots, x_t) \in \mathbb{R}^t$ is measurable, and so forth). The probability measures on the product space $\mathcal{X} \times \mathcal{Y}$ with marginals $\mu, \nu$ correspond to all possible transport plans between the given marginals. Denote this set

$$\Pi(\mu, \nu) = \{\gamma \in \mathcal{P}(\mathcal{X} \times \mathcal{Y}) \text{ with marginals } \mu \text{ and } \nu\}.$$

We will often consider pairs of random variables $(X, Y)$ defined on some probability space $(\Omega, \mathbb{P})$ and taking values in resp. $\mathcal{X}$ and $\mathcal{Y}$, and refer to them as transport plans as well. Any property on $(X, Y)$ should then be understood as a property on $(X, Y)_* \mathbb{P}$.

In the following we assume w.l.o.g. that $\mathcal{X} = supp(\mu)$ and $\mathcal{Y} = supp(\nu)$, whenever dealing with transport problems between $\mu$ and $\nu$. With some abuse of notation we will often write $\mathbb{R}^N$, the reader keeping in mind we mean $\mathcal{X}$ or $\mathcal{Y}$.[1] For the sake of simplicity, for us being measurable with respect to a sigma algebra means to be equal to a correspondingly measurable function modulo a null set w.r.t. the measure unequivocally relevant to the given context.

**Definition 2.1.** *A transport plan $\gamma \in \Pi(\mu, \nu)$ is called[2] causal (between $\mu$ and $\nu$) if for any $t \in \mathbb{T}$ and $B \in \mathcal{F}_t^{\mathcal{Y}}$, the mapping $x \in \mathcal{X} \to \gamma^x(B)$ is $\mathcal{F}_t^{\mathcal{X}}$-measurable. The set of all such plans will be denoted*

$$\Pi_c(\mu, \nu).$$

Analogously, we will be interested in transport plans that are "causal in both directions", or bicausal in our terminology. The set of all such plans is explicitly given by

$$\Pi_{bc}(\mu, \nu) = \{\gamma \in \Pi_c(\mu, \nu) \text{ s.t. } e_*\gamma \in \Pi_c(\nu, \mu)\},$$

where $e(x, y) = (y, x)$. As in the usual optimal transport problem, the set of all causal plans $\Pi_c(\mu, \nu)$, as well as $\Pi_{bc}(\mu, \nu)$, are always non-empty because $\mu \otimes \nu \in \Pi_{bc}(\mu, \nu)$. Further, as in the classical setting, we shall consider the case of Borel measurable transformations $T : \mathcal{X} \to \mathcal{Y}$ satisfying $T_*\mu = \nu$, so that in particular $\gamma^T := (id \times T)_*\mu$ belongs to $\Pi(\mu, \nu)$, and call them (Monge) transport maps. Transport maps are termed (bi)causal if the associated $\gamma^T$ is so.

The condition for a transport plan to be causal, written in terms of stochastic processes, looks as follows: for all $t = 1, \ldots, N$ and $B_t \in \mathcal{F}_t^{\mathcal{Y}}$,

$$\mathbb{P}(Y_1 \in B_1, \ldots, Y_t \in B_t \mid X_1, \ldots, X_N) = \mathbb{P}(Y_1 \in B_1, \ldots, Y_t \in B_t \mid X_1, \ldots X_t).$$

---

[1] Throughout this work most of the results would still hold for $S$-valued discrete-time stochastic processes in $N$-steps, with $S$ a Polish space. That is, we could take $\mathcal{X} = \mathcal{Y} = S^N$. For notational simplicity we take $S = \mathbb{R}$.

[2] Lassalle ([Las15]) introduced more technical definitions. The one here is enough for us.

Heuristically this reads as "given the past of $X$, the past of $Y$ and the future of $X$ are independent". This is perhaps best interpreted by the following equivalent formulation (see e.g. [Kal02, Proposition 6.13]):

$$Y_t = F_t(X_1, \ldots, X_t, U_t) \quad \forall t \in \{1, \ldots, N\},$$

for some measurable functions $F_t$ and where each $U_t$ is a uniform random variable independent of $X_1, \ldots, X_N$.

**Remark 2.2.** *Clearly a transport map $T$ is causal if and only if it is adapted, in the sense that there exist Borel-measurable $T^t : \mathbb{R}^t \to \mathbb{R}$ such that for $\mu$-a.e. $(x_1, \ldots, x_N)$:*

$$T(x_1, \ldots, x_N) \;=\; (T^1(x_1), T^2(x_1, x_2), \ldots, T^N(x_1, \ldots, x_N)).$$

The following proposition allows us on the one hand to characterize the causal transport plans using the successive disintegrations of measures on a product space. On the other hand, it shows that causality can be seen as a linear constraint on measures on the product space, stated in terms of a special class of test functions or via discrete stochastic integrals.

**Proposition 2.3.** *The following statements are equivalent:*

  *(1) $\gamma$ is a causal transport plan on $\mathcal{X} \times \mathcal{Y}$ between the measures $\mu$ and $\nu$.*
  *(2) Decomposing $\gamma$ in terms of successive regular kernels*

$$(2.1) \qquad \gamma(dx_1, \ldots, dx_N, dy_1, \ldots, dy_N) = \bar{\gamma}(dx_1, dy_1)\gamma^{x_1, y_1}(dx_2, dy_2) \ldots \gamma^{x_1, \ldots, x_{N-1}, y_1, \ldots, y_{N-1}}(dx_N, dy_N),$$

  *then $\bar{\gamma} \in \Pi(p_*^1 \mu, p_*^1 \nu)$ and for $t < N$ and $\gamma$-almost all $x_1, \ldots, x_t, y_1, \ldots, y_t$*

$$(2.2) \qquad\qquad p_*^1 \gamma^{x_1, \ldots, x_t, y_1, \ldots, y_t} \;=\; \mu^{x_1, \ldots, x_t},$$

  *and for $\nu$-almost all $y_1, \ldots, y_t$*

$$(2.3) \qquad\qquad \gamma^{y_1, \ldots, y_t}(dy_{t+1}) = \nu^{y_1, \ldots, y_t}(dy_{t+1}).$$

  *(3) $\gamma \in \Pi(\mu, \nu)$ and for all $t \in \{1, \ldots, N\}$, $h_t \in C_b(\mathbb{R}^t)$ and $g_t \in C_b(\mathbb{R}^N)$ we have*

$$\int h_t(y_1, \ldots, y_t) \left\{ g_t(x_1, \ldots, x_N) - \int g_t(x_1, \ldots, x_t, \bar{x}_{t+1}, \ldots, \bar{x}_N)\mu^{x_1, \ldots, x_t}(d\bar{x}_{t+1}, \ldots, d\bar{x}_N) \right\} d\gamma = 0.$$

  *(4) $\gamma \in \Pi(\mu, \nu)$ and for every bounded continuous $\mathcal{F}^{\mathcal{Y}}$-adapted process $H$ and each bounded $\mathcal{F}^{\mathcal{X}}$-martingale $M$ we have*

$$\int \sum_{t < N} H_t(y_1, \ldots, y_t) \left[ M_{t+1}(x_1, \ldots, x_{t+1}) - M_t(x_1, \ldots, x_t) \right] d\gamma = 0.$$

The proof of the above result is given in the next section. Notice that (2.3) is equivalent to

$$\int \gamma^{x_1, \ldots, x_t, y_1, \ldots, y_t}(\mathbb{R}, dy_{t+1})\gamma^{y_1, \ldots, y_t}(dx_1, \ldots, dx_t) \;=\; \nu^{y_1, \ldots, y_t}(dy_{t+1}),$$

which is more convenient for the derivation of the dynamic programming principle to come.

We now introduce our main optimization problem, the causal optimal transport problem: given some Borel cost function $c$ defined on $\mathbb{R}^N \times \mathbb{R}^N$ and the probability measures $\mu, \nu$, find the minimal cost at which they can be coupled in a causal way, i.e. consider

$$(\text{Pc}) \qquad\qquad \inf_{\gamma \in \Pi_c(\mu, \nu)} \int c \, d\gamma.$$

Minimizing over the set $\Pi_{bc}(\mu, \nu)$ defines the bicausal optimal transport problem

$$(\text{Pbc}) \qquad\qquad \inf_{\gamma \in \Pi_{bc}(\mu, \nu)} \int c \, d\gamma.$$

These should be compared to the classical problem of optimal transport in which minimization is done over $\Pi(\mu, \nu)$. Let us introduce an assumption, which eases the proof of Theorem 2.5:

**Assumption 2.4.** *The measure $\mu$ is successively weakly continuous in the sense that for each $t < N$, there is a version of the regular conditional kernel of $\mu$ w.r.t. its first $t$ variables s.t.*

$$(x_1, \ldots, x_t) \in supp(\mu) \cap \mathbb{R}^t \mapsto \mu^{x_1, \ldots, x_t}(d\bar{x}_{t+1}, \ldots, d\bar{x}_N) \in \mathcal{P}(\mathbb{R}^{N-t}),$$

*is continuous w.r.t. the weak topology in the range and the relative topology in the domain.*

Let us observe that if $supp(\mu)$ contains no accumulation points, as in the random walk/event tree setting, Assumption 2.4 is vacuously fulfilled. On the other extreme, there are many discrete-time processes with full support satisfying it, e.g. the Gaussian case. Let us also notice that $\mu^{x_1,\dots,x_{N-1}}$ is automatically a univariate measure, and more generally we will commonly write $\mu^{x_1,\dots,x_t}(d\bar{x}_{t+1},\dots,d\bar{x}_{t+k})$ for the measure $\mu^{x_1,\dots,x_t}(d\bar{x}_{t+1},\dots,d\bar{x}_{t+k},\mathbb{R}^{N-t-k})$ and similarly $\mu(dx_1,\dots,dx_t)$ for the projection of $\mu$ into the first $t$-marginals. The following sets of test functions will be instrumental for the dual formulation:

$$(2.4) \qquad \mathbb{F} := \left\{ \begin{array}{c} F : \mathbb{R}^N \times \mathbb{R}^N \to \mathbb{R} \text{ s.t. } F(x_1,\dots,x_N,y_1,\dots,y_N) = \\ \sum_{t<N} h_t(y_1,\dots,y_t)\left[g_t(x_1,\dots,x_N) - \int g_t(x_1,\dots,x_t,\bar{x}_{t+1},\dots,\bar{x}_N)\mu^{x_1,\dots,x_t}(d\bar{x}_{t+1},\dots,d\bar{x}_N)\right], \\ \text{with } h_t \in C_b(\mathbb{R}^t), g_t \in C_b(\mathbb{R}^N) \text{ for all } t < N \end{array} \right\},$$

$$(2.5) \qquad \mathbb{S} := \left\{ \begin{array}{c} S : \mathbb{R}^N \times \mathbb{R}^N \to \mathbb{R} \text{ s.t. } S(x_1,\dots,x_N,y_1,\dots,y_N) = \\ \sum_{t<N} H_t(y_1,\dots,y_t)\left[M_{t+1}(x_1,\dots,x_{t+1}) - M_t(x_1,\dots,x_t)\right], \\ \text{with } H_t, M_t \in C_b(\mathbb{R}^t) \text{ for all } t < N, \text{ and with } M \text{ a martingale} \end{array} \right\}.$$

All in all we are ready to present the basic primal attainability/no-duality-gap result:

**Theorem 2.5.** *Suppose that $c : \mathbb{R}^N \times \mathbb{R}^N \to \mathbb{R} \cup \{+\infty\}$ is lower semicontinuous and bounded from below, and that Assumption 2.4 holds. Then there is no duality gap*

$$\inf_{\gamma \in \Pi_c(\mu,\nu)} \int c\,d\gamma = \sup_{\substack{\Phi,\Psi \in C_b(\mathbb{R}^N), F \in \mathbb{F} \\ \Phi \oplus \Psi \leq c + F}} \left[\int \Phi\,d\mu + \int \Psi\,d\nu\right] = \sup_{\substack{\Phi,\Psi \in C_b(\mathbb{R}^N), S \in \mathbb{S} \\ \Phi \oplus \Psi \leq c + S}} \left[\int \Phi\,d\mu + \int \Psi\,d\nu\right],$$

*and the infimum on the l.h.s. (i.e. (Pc)) is attained.*

We observe that Assumption 2.4 allows us to test against continuous bounded functions, instead of just bounded Borel, which is necessary for the simple proof of the previous theorem given in Section 3. However, this assumption can be lifted at the price of losing simplicity (see [Las15], which was written concurrently), and we choose not to prove the most general statement as our interest lies in other aspects of the problem. It is also easy to see that the dual problem can be reduced to both:

$$\sup_{\Psi \in C_b(\mathbb{R}^N), F \in \mathbb{F}, \Psi \leq c + F} \int \Psi\,d\nu \quad \text{and} \quad \sup_{\Psi \in C_b(\mathbb{R}^N), S \in \mathbb{S}, \Psi \leq c + S} \int \Psi\,d\nu.$$

The analogue of this theorem for bicausal transport plans is given in the next section, and was first obtained in [PP12, Theorem 7.2]

Going back to Proposition 2.3, the importance of decomposition (2.1) lies in the fact that it suggests that the causal optimal transport problem can be solved recursively if the starting measure $\mu$ is Markovian and the cost function has a "semiseparable" structure:

**Theorem 2.6** (DPP for causal plans)**.** *Let $\mu, \nu \in \mathcal{P}(\mathbb{R}^N)$, suppose that $\mu$ is a Markov measure and that the cost is semiseparable in the sense that for non-negative l.s.c. functions $c_t$ we have $c = \sum_t c_t(x_t, y_1,\dots,y_t)$. Then, starting from $V_N^c := 0$ and defining recursively for $t = N,\dots,2$:*

$$(2.6) \quad V_{t-1}^c(y_1,\dots,y_{t-1}; m(dx_{t-1})) :=$$

$$\inf_{\gamma \in \Pi\left(\int_{x_{t-1}} m(dx_{t-1})\mu^{x_{t-1}}(dx_t), \nu^{y_1,\dots,y_{t-1}}(dy_t)\right)} \int \gamma(dx_t, dy_t)\Big\{c_t(x_t, y_1,\dots,y_t)$$

$$+ V_t^c(y_1,\dots,y_t; \gamma^{y_t}(dx_t))\Big\},$$

*and*

$$V_0^c := \inf_{\gamma \in \Pi(p_*^1\mu, p_*^1\nu)} \int \gamma(dx_1, dy_1)\{c_1(x_1, y_1) + V_1^c(y_1; \gamma^{y_1}(dx_1))\},$$

*the integrals above are well-defined and $V_0^c = value(\text{Pc})$, i.e. the recursion determines (Pc). Furthermore, if Assumption 2.4 holds, there exists a causal optimizer $\tilde{\gamma}$ with the additional property that for all $t$:*

$$\tilde{\gamma}^{x_1,\dots,x_t,y_1,\dots,y_t}(dx_{t+1}, dy_{t+1}) = \tilde{\gamma}^{x_t,y_1,\dots,y_t}(dx_{t+1}, dy_{t+1}),$$

*and each of the one-step optimization problems in (2.6) corresponds to a general transport problem of [GRST15] which is convex and l.s.c. in the kernel (i.e. in $\gamma^{y_t}$, for fixed $y_1,\dots,y_{t-1}$).*

The very last line of the previous result means concretely that

$$(2.7) \quad V_{t-1}^c(y_1, \ldots, y_{t-1}; m(dx_{t-1})) =$$

$$\inf_{\substack{y_t \mapsto \gamma^{y_t} \ s.t. \\ \int_{y_t} \nu^{y_1, \ldots, y_{t-1}}(dy_t)\gamma^{y_t}(dx_t) = \int_{x_{t-1}} m(dx_{t-1})\mu^{x_{t-1}}(dx_t)}} \int_{y_t} \nu^{y_1, \ldots, y_{t-1}}(dy_t) \Big\{ \int_{x_t} \gamma^{y_t}(dx_t)c_t(x_t, y_1, \ldots, y_t)$$

$$+ V_t^c(y_1, \ldots, y_{t-1}, y_t; \gamma^{y_t}(dx_t)) \Big\},$$

and that the function in curly brackets in the r.h.s. is convex and l.s.c. in $\gamma^{y_t}$ (ceteris paribus). The same function in brackets is only jointly universally measurable (actually lower semianalytic) as far as we can say, regardless of Assumption 2.4. If on the other hand Assumption 3.2 holds, then this function is jointly l.s.c. In [GRST15] problems of this kind have been analysed and their duality theory has been established; we provide the dual formulation of $V_{t-1}^c$ in (4.7).

Our next main result, Theorem 2.7 (which we prove in Section 5), establishes the equivalence of (Pc) and (Pbc) under Condition 2.11 below, and furthermore, it identifies the causal optimizer of (Pc) for convex separable costs (as well as clarifying when these are Monge maps), in a setting relevant for future applications.

Let us introduce some useful notation. By $F_\eta$ we mean the usual distribution function of a probability measure $\eta$ on the line (we denote $F_{\nu_1}$ the distribution of $p_*^1\nu$), and by $F_\eta^{-1}(u)$ its left-continuous generalized inverse, i.e. $F_\eta^{-1}(u) = \inf\{y : F_\eta(y) \geq u\}$. Also the increasing $N$-dimensional **Knothe-Rosenblatt** rearrangement of $\mu$ and $\nu$ is defined as the law of the random vector $(X_1^*, \ldots, Y_N^*, X_1^*, \ldots, Y_N^*)$ where

$$(2.8) \qquad X_1^* = F_{\mu_1}^{-1}(U_1), \qquad Y_1^* = F_{\nu_1}^{-1}(U_1), \qquad \text{and inductively}$$
$$X_n^* = F_{\mu^{X_1^*, \ldots, X_{n-1}^*}}^{-1}(U_n), \quad Y_n^* = F_{\nu^{Y_1^*, \ldots, Y_{n-1}^*}}^{-1}(U_n), \text{ for } n = 2, \ldots, N,$$

for $U_1, \ldots, U_N$ independent and uniformly distributed random variables on $[0, 1]$. Additionally, if $\mu$-a.s. all the conditional distributions of $\mu$ are atomless (e.g. if $\mu$ has a density), then this rearrangement is induced by the (Monge) map

$$(x_1, \ldots, x_N) \mapsto T(x_1, \ldots, x_N) := (T^1(x_1), \ldots, T^N(x_N; x_1, \ldots x_{N-1})),$$

where

$$T^1(x_1) := F_{\nu_1}^{-1} \circ F_{\mu_1}(x_1),$$
$$(2.9) \qquad T^n(x_n; x_1, \ldots, x_{n-1}) := F_{\nu^{T^1(x_1), \ldots, T^{n-1}(x_{n-1}; x_1, \ldots, x_{n-2})}}^{-1} \circ F_{\mu^{x_1, \ldots, x_{n-1}}}(x_n), \quad n \geq 2.$$

**Theorem 2.7.** *Assume that $c$ is l.s.c. bounded from below and has a separable structure*

$$(2.10) \qquad c(x_1, \ldots, x_N, y_1, \ldots, y_N) = \sum_{t \leq N} c_t(x_t, y_t).$$

*Further suppose that the starting measure $\mu$ is the product of its marginals, i.e.*

$$(2.11) \qquad \mu(dx_1, \ldots, dx_N) = \mu_1(dx_1) \ldots \mu_N(dx_N).$$

*Then the values of (Pc) and (Pbc) coincide. If moreover, $c_t(x, y) = c_t(x - y)$ and $c_t$ is convex[3], then a solution to (Pc) is given by the Knothe-Rosenblatt rearrangement Additionally, if each $\mu_i$ is atomless (e.g. if they have a density), then this rearrangement is induced by the Monge map (2.9).*

**Corollary 2.8.** *Assume that $c$ is l.s.c. bounded from below and is of the form*

$$(2.12) \qquad c(x_1, \ldots, x_N, y_1, \ldots, y_N) = c_1(x_1, y_1) + \sum_{1 < t \leq N} c_t(x_t - x_{t-1}, y_t - y_{t-1}),$$

*and the source measure $\mu$ has independent increments. Then the values of (Pc) and (Pbc) coincide. If moreover, $c_t(a, b) = c_t(a - b)$ and $c_t$ is convex then a solution to (Pc) is given by the Knothe-Rosenblatt rearrangement (2.8) and if additionally each $\mu_i$ is atomless then this rearrangement is induced by the Monge map (2.9).*

---

[3]More generally, a Spence-Mirrlees condition suffices.

In [Las13, Lemma 5] it was shown that the optimal solution to the causal transport problem in Wiener Space, with Cameron-Martin cost and Wiener measure as source, is bicausal. The results above give the causal/bicausal equality and existence as well as characterization of the Monge solution to what can be thought of as "finite-dimensional projections" of that problem. We can only guess then that causal/bicausal equality in continuous-time is a result of the Cameron-Martin cost being written in terms of "speed" and the source measure having independent increments. In Section 7 we give two counterexamples showing that if either independence or separability of the cost is dropped, the causal-bicausal equality may fail.

We introduce now our last main results. The goal here is to explore and to further the connection between (non-anticipative) multistage stochastic programming and bicausal optimal transport, discovered by Pflug and Pichler [PP12], from the point of view of geometric/functional inequalities. We start by defining the value function of a stochastic program

$$(2.13) \quad v(\eta) := \inf_{u_1(),\dots,u_N()} \int H(x_1,\dots,x_N, u_1(x_1), u_2(x_1,x_2),\dots, u_N(x_1,\dots,x_N))\eta(dx_1,\dots,dx_N),$$

as a function of the noise distribution $\eta$. Here $H : \mathbb{R}^N \times \mathbb{R}^N \to \mathbb{R}$ is the objective function and the minimization is taken over Borel adapted controls. As established in [PP12, Theorem 6.1], as soon as $H$ is 1-Lipschitz in its first argument and convex in the second one, the discrepancy $|v(\mu) - v(\nu)|$ is less than the bicausal distance between $\mu$ and $\nu$ w.r.t. the cost

$$c(x,y) := \|x - y\|_1 = \sum_i |x_i - y_i|.$$

This means that if say $\mu$ is a "complicated law" (for instance given by a nightmarish tree), then finding a "simpler law" $\nu$ close enough in the causal distance sense guarantees that uniformly in the given class of stochastic programs the discrepancy is controlled. Therefore, the practical question is how to efficiently minimize this bicausal distance over a given set of measures. We argue that in many situations a transport-information inequality for such bicausal distances permits to gauge the discrepancy of stochastic programs replacing the computation of such distance by the evaluation of a relative entropy. Notice that in our interpretation, the complicated measure $\mu$ is expected to dominate (in the sense of absolute continuity) the simpler measure $\nu$. Our result is reminiscent of [PP14, Proposition 4.6]; crucially the presence of the entropy here stems from our assumption (EXP).

**Theorem 2.9.** *Let $\mu \in \mathcal{P}(\mathbb{R}^N)$ satisfy:*

(EXP) *for each $t$ there exist $a_t > 0, \lambda_t \in \mathbb{R}$ such that $\int e^{a_t x_t^2} \mu^{x_1,\dots,x_{t-1}}(dx_t) \leq e^{\lambda_t}$ $\mu$-a.s.*

(LIP) *There is $C > 0$ such that for all $t \in \{1,\dots,N\}$, and every 1-Lipschitz function $f : \mathbb{R} \to \mathbb{R}$, the function*

$$(x_1,\dots,x_{t-1}) \mapsto \int f(x_t)\mu^{x_1,\dots,x_{t-1}}(dx_t),$$

*is $\mu$-a.s. equal to a $C$-Lipschitz function.*

*Then we have the following bicausal transport-information inequality:*

$$(2.14) \qquad \text{for all } \nu \in \mathcal{P}(\mathbb{R}^N): \quad \mathcal{W}_{1,bc}(\mu,\nu) := \inf_{\gamma \in \Pi_{bc}(\mu,\nu)} \int \|x - y\|_1 d\gamma(x,y) \leq K\sqrt{Ent(\nu|\mu)},$$

*where $Ent(\cdot|\mu)$ denotes the relative entropy with respect to $\mu$ and*

$$K := \sqrt{2 \sum_{j < N} (1 + C)^{2j} \frac{(1 + \lambda_{N-j})}{a_{N-j}^2}}.$$

*In particular, for every "cost criterion" $H : \mathbb{R}^N \times \mathbb{R}^N \to \mathbb{R}$ bounded from below, $r$-Lipschitz in its first argument and convex in its second, we have for the corresponding stochastic programs*

$$(2.15) \qquad\qquad\qquad |v(\mu) - v(\nu)| \leq rK\sqrt{Ent(\nu|\mu)},$$

*with $v(\cdot)$ defined as in (2.13).*

**Remark 2.10.** *A few observations are in order:*

(1) *W.l.o.g. the Lipschitz property above is meant with respect to the sum-of-absolute-values norms in the respective spaces. By the tower property $(EXP)$ implies*

$$\int e^{a_1 x_1^2 + \dots + a_N x_N^2} \mu(dx_1,\dots,dx_N) \leq e^{\lambda_1 + \dots + \lambda_N} < \infty,$$

*and in particular every Lipschitz function is $\mu$-integrable.*

(2) *We stress that Assumption* $(EXP)$ *is reasonable in practice; it is automatically satisfied in the discrete case, for empirical measures, or when* $\mu$ *has bounded support. The corresponding non-causal version of* (2.14), *which can be found in* [Vil08, Theorem 22.10] *or* [DGW04, Theorem 2.3] *and crucially does not implies ours, actually also needs a condition equivalent to* $(EXP)$.

(3) *Assumption* $(LIP)$ *implies, by the Kantorovich-Rubinstein Theorem (e.g.* [Vil03, Theorem 1.14]) *that*

$$\mathcal{W}_1(\mu^{x_1,\dots,x_{t-1}}, \mu^{z_1,\dots,z_{t-1}}) \leq C \sum_{i<t} |x_i - z_i|,$$

*where* $\mathcal{W}_1$ *is the usual 1-Wasserstein distance on* $\mathcal{P}(R)$. *If* $\mu$ *is a Markov law, then the r.h.s. here can be taken as* $C|x_{t-1} - z_{t-1}|$. *On the other hand, if* $\mu$ *is a martingale law, the l.h.s. is in turn bounded from below by* $|x_{t-1} - z_{t-1}|$.

(4) *All in all, Assumption* $(LIP)$ *is the most fundamental for Theorem 2.9. It is however clearly satisfied in the discrete case or for empirical measures. In general it is also implied if* $\mu$ *has independent increments (with* $C = 1$ *then), or if e.g.* $\mu$ *is the law of the solution* $X_1, \dots, X_N$ *of a uniformly Lipschitz random dynamical system of the form* $X_{t+1} = R(Z_t, (X_1, \dots, X_t), t)$, *where* $R$ *is Borel in the first argument, Lipschitz in the second one (uniformly w.r.t. the first one), and* $Z_t$ *is independent of* $(X_1, \dots, X_t)$.

The proof of the previous theorem is given in Section 6. We observe that the constant in (2.14) can be very plausibly improved (making it of order $\sqrt{N}$ is probably the best possible). This is done for example in [DGW04, Theorem 2.5] for the non-causal but Markov case, and necessitates more care and further assumptions.

In Section 6 we will also connect other modern tenant of optimal transport theory with the field of stochastic programming, and more broadly, nested distances and the bicausal setting. Inspired by the concepts of displacement interpolation/convexity (as in [Vil03, Chapter 5]), we will define a related notion where the Knothe-Rosenblatt rearrangement replaces the role of the Brenier's map, and dub it lexicographical displacement interpolation/convexity. We expect that the lexicographical displacement interpolations should have a geometric meaning as geodesic curves, however the materialization of this idea is left as an open problem, which we consider as a relevant step towards an understanding of the geometrical side of the bicausal/nested transport problem. This aside, we concretely show that under convexity and Lipschitz conditions on the cost, stochastic programs are lexicographical displacement concave in analogy to the potential energy in optimal transport ([Vil03, Chapter 5.2]).

## 3. Duality

The dual to the causal transport problem is discussed in this section. First, we give the proof of Proposition 2.3.

*Proof of Proposition 2.3.* STEP 1: Equivalence between Points 1 and 3:

Denote $f^h(x_1, \dots, x_N) := \int h_t(y_1, \dots, y_t) \gamma^{x_1, \dots, x_N}(dy_1, \dots, dy_t)$ with $h_t \in C(\mathbb{R}^t)$. By definition $\gamma \in \Pi_c(\mu, \nu)$ if and only if for all $t \leq N$ and all such $f^h$ we have

$$f^h(x_1, \dots, x_N) = \int f^h(x_1, \dots, x_t, \bar{x}_{t+1}, \dots, \bar{x}_N) \mu^{x_1, \dots, x_t}(d\bar{x}_{t_1}, \dots, d\bar{x}_N),$$

which is equivalent to the following:

$$\int g(x_1, \dots, x_N) \left[ f^h(x_1, \dots, x_N) - \int f^h(x_1, \dots, x_t, \bar{x}_{t+1}, \dots, \bar{x}_N) \mu^{x_1, \dots, x_t}(d\bar{x}_{t_1}, \dots, d\bar{x}_N) \right] d\mu = 0,$$

for every function $g \in C_b(\mathbb{R}^N)$ and for all $t \leq N$. The fact we can take the $g$'s continuous and not merely Borel bounded comes from the fact that $\mu$ is a Borel finite measure on a Polish space. It is easy to see that the previous equation is equivalent to

$$\int f^h(x_1, \dots, x_N) \left[ g(x_1, \dots, x_N) - \int g(x_1, \dots, x_t, \bar{x}_{t+1}, \dots, \bar{x}_N) \mu^{x_1, \dots, x_t}(d\bar{x}_{t_1}, \dots, d\bar{x}_N) \right] d\mu = 0.$$

Finally, by the tower property of conditional expectations the latter is equivalent to:

$$\int h_t(y_1, \dots, y_t) \left[ g_t(x_1, \dots, x_N) - \int g_t(x_1, \dots, x_t, \bar{x}_{t+1}, \dots, \bar{x}_N) \mu^{x_1, \dots, x_t}(d\bar{x}_{t+1}, \dots, d\bar{x}_N) \right] d\gamma \quad = \quad 0.$$

STEP 2: Equivalence between Points 1 and 3:

Let $\gamma \in \Pi(\mu, \nu)$ be decomposed as in (2.1). By induction, one readily sees that causality is equivalent to the conditions (for all $t$) Thus for the equivalence of points 1. and 2. we only need to argue that

knowing $\gamma \in \Pi(\mu, \nu)$ implies (2.3) and then conversely, that knowing (2.3), (2.2) and $\bar{\gamma} \in \Pi(p_*^1 \mu, p_*^1 \nu)$ implies $\gamma \in \Pi(\mu, \nu)$. The direct statement is clear, since in the l.h.s. of (2.3) the dependence on the $x$'s is integrated out, so what is left is $\gamma^{y_1, \ldots, y_t}(dy_{t+1})$, which must coincide with a kernel of $\nu$. For the converse statement, the fact that the $x$-marginal of $\gamma$ is $\mu$ follows from (2.2) and $\bar{\gamma} \in \Pi(p_*^1 \mu, p_*^1 \nu)$, whereas the fact that the $y$-marginal of $\gamma$ is $\nu$ follows from $\bar{\gamma} \in \Pi(p_*^1 \mu, p_*^1 \nu)$ and (2.2) after the above observation on the l.h.s. of such equation.

STEP 3: Equivalence between Points 3 and 4:

Evidently in Point 3 we could have taken $h_t$ and $g_t$ Borel bounded, as STEP 1 suggests. Choosing then $g_t = M_{t+1}$ and $h_t = H_t$ for each $t < N$, and summing up, proves Point 4 from Point 3. Conversely, given $h_t$ and $g_t$ we build $H_s = h_t \mathbf{1}_{\mathbf{s} \geq \mathbf{t}}$ and $M_s = \int g_t(x_1, \ldots, x_s, x_{s+1}, \ldots, x_N) \mu^{x_1, \ldots, x_s}(dx_{s+1}, \ldots, dx_N)$ and conclude by telescopic sum. $\qquad \square$

We now proceed to the duality and attainment questions. First:

**Lemma 3.1.** *Suppose Assumption 2.4 holds. For $g \in C_b(\mathbb{R}^N)$, the functions*

$$(x_1, \ldots, x_N) \mapsto g(x_1, \ldots, x_N) - \int g(x_1, \ldots, x_t, \bar{x}_{t+1}, \ldots, \bar{x}_N) \mu^{x_1, \ldots, x_t}(d\bar{x}_{t+1}, \ldots, d\bar{x}_N),$$

*belong themselves to $C_b(\mathbb{R}^N)$.*

*Proof.* It is suffices to check the continuity of

$$(x_1, \ldots, x_t) \mapsto \int g(x_1, \ldots, x_t, \bar{x}_{t+1}, \ldots, \bar{x}_N) \mu^{x_1, \ldots, x_t}(d\bar{x}_{t+1}, \ldots, d\bar{x}_N).$$

Let $x^n := (x_1^n, \ldots, x_t^n)$ converge to $y := (y_1, \ldots, y_t)$, thus we may assume all $x^n$ belongs to a fixed ball $B$ around $y^n$. Let us fix an arbitrary $\varepsilon > 0$. By assumption $\mu^{x^n}$ converges weakly to $\mu^y$ and so in particular $\{\mu^{x^n}\}_{n \in \mathbb{N}}$ is tight. Thus we may find a compact in $\mathbb{R}^{N-t}$ such that $\sup_n \mu^{x^n}(K^c) \leq \varepsilon/3$. For large enough $n_0$ we also have that $\sup_{n > n_0, z \in K} |g(x^n, z) - g(y, z)| \leq \varepsilon/3$, since $g$ restricted to $B \times K$ must be uniformly continuous. We then write:

$$\left| \int g(x^n, z) \mu^{x^n}(dz) - \int g(y, z) \mu^y(dz) \right| \leq \alpha + \beta + \gamma,$$

with

$$\alpha = \left| \int_K [g(x^n, z) - g(y, z)] \mu^{x^n}(dz) \right|, \beta = \left| \int_{K^c} [g(x^n, z) - g(y, z)] \mu^{x^n}(dz) \right|, \gamma = \left| \int g(y, z) [\mu^{x^n}(dz) - \mu^y(dz)] \right|.$$

We easily see that each term is smaller that $\varepsilon/3$ for $n$ large enough. $\qquad \square$

*Proof of Theorem 2.5.* By Proposition 2.3, item 3, we know that the set $\Pi_c(\mu, \nu)$ is the intersection of the compact $\Pi(\mu, \nu)$ with all the subspaces defined by the functions in (2.4), each of them closed owing to Lemma 3.1. Thus $\Pi_c(\mu, \nu)$ is compact and primal attainment follows easily. If the cost function belongs to $C_b(\mathbb{R}^N \times \mathbb{R}^N)$, this theorem can be seen as a particular case of the Monge-Kantorovich problem with additional linear constraints considered in [Zae15, Theorem 2.5], where we just have to show that $\mathbb{F} \subset C_b(\mathbb{R}^N \times \mathbb{R}^N)$; but this is clear again by Lemma 3.1. On the other hand, when working with $\mathbb{S} \subset C_b(\mathbb{R}^N \times \mathbb{R}^N)$, it is easy to see that the same lemma implies that the martingales $M$ can be assumed continuous in the context of Proposition 2.3 and its proof, so again [Zae15, Theorem 2.5] establishes our stated result. Finally duality for lower semicontinuous cost functions is achieved by the usual approximation arguments in [Vil03, Chapter 1], owing to the compactness of the set of all causal plans. $\qquad \square$

The analogue of Theorem 2.5 is obtained for the bicausal setting. As in the causal case, we define the set

$$(3.1) \qquad \mathbb{F}' := \left\{ \begin{array}{c} F : \mathbb{R}^N \times \mathbb{R}^N \to \mathbb{R} \text{ s.t. } F(x_1, \ldots, x_N, y_1, \ldots, y_N) = \\ \sum_{t < N} h_t(y_1, \ldots, y_t) \left[ g_t(x_1, \ldots, x_N) - \int g_t(x_1, \ldots, x_t, \bar{x}_{t+1}, \ldots, \bar{x}_N) \mu^{x_1, \ldots, x_t}(d\bar{x}_{t+1}, \ldots, d\bar{x}_N) \right] + \\ \sum_{t < N} h_t'(x_1, \ldots, x_t) \left[ g_t'(y_1, \ldots, y_N) - \int g_t'(y_1, \ldots, y_t, \bar{y}_{t+1}, \ldots, \bar{y}_N) \nu^{y_1, \ldots, y_t}(d\bar{y}_{t+1}, \ldots, d\bar{y}_N) \right], \\ \text{with } h_t, h_t' \in C_b(\mathbb{R}^t), \ g_t, g_t' \in C_b(\mathbb{R}^N) \text{ for all } t < N \end{array} \right\}.$$

The following strengthened version of Assumption 2.4 will be needed:

**Assumption 3.2.** *Both $\mu$ and $\nu$ are successively weakly continuous.*

**Corollary 3.3.** *Suppose that* $c : \mathbb{R}^N \times \mathbb{R}^N \to \mathbb{R}$ *is lower semicontinuous and bounded from below, and that Assumption 3.2 holds. Then there is no duality gap:*

(Dbc)
$$\inf_{\gamma \in \Pi_{bc}(\mu,\nu)} \int c d\gamma = \sup_{\substack{\Phi,\Psi \in C_b(\mathbb{R}^N), F' \in \mathbb{F}' \\ \Phi \oplus \Psi \leq c + F'}} \left[ \int \Phi d\mu + \int \Psi d\nu \right].$$

*Moreover, the infimum in the l.h.s. is attained.*

We omit the obvious formulation with discrete stochastic integrals, as well as the proof.

**Remark 3.4.** *We have observed that a measurable* $T : \mathbb{R}^N \to \mathbb{R}^N$ *is causal if and only if it is adapted, in the sense that* $\mu$*-a.s.* $T(x_1, \ldots, x_N) = (T^1(x_1), T^2(x_1, x_2), \ldots, T^N(x_1, \ldots, x_N))$ *for measurable* $T^i :$ $\mathbb{R}^t \to \mathbb{R}$. *Bicausality is more subtle, but clearly holds if for instance* $T$ *admits a measurable* $\mu$*-a.s. left inverse which is also adapted.*

## 4. DYNAMIC PROGRAMMING PRINCIPLE

Let $\gamma \in \Pi(\mathbb{R}^N, \mathbb{R}^N)$. It is then possible to decompose $\gamma$, uniquely in a suitable way, by successive disintegration first w.r.t. $(x_1, \ldots, x_{N-1}, y_1, \ldots, y_{N-1})$, then w.r.t. $(x_1, \ldots, x_{N-2}, y_1, \ldots, y_{N-2})$ and so forth in a recursive way, obtaining

(4.1) $\quad \gamma(dx_1, \ldots, dx_N, dy_1, \ldots, dy_N) = \bar{\gamma}(dx_1, dy_1)\gamma^{x_1, y_1}(dx_2, dy_2) \ldots \gamma^{x_1, \ldots, x_{N-1}, y_1, \ldots, y_{N-1}}(dx_N, dy_N),$

so that each $\gamma^{x_1, \ldots, x_t, y_1, \ldots, y_t}(dx_{t+1}, dy_{t+1})$ is a regular conditional probability or kernel (see [Bog07, Theorem 10.4.14, Corollary 10.4.17] and [Kal02, Lemma 1.41 ]). We will freely perform concatenation of these objects or further decompose them into smaller kernels, as justified in [BS78, Chapter 7.4.3]. The characterization of causal plans through the kernels was given in Proposition 2.3.

We believe that generally (Pc) does not allow for a meaningful and useful "factorized" or "recursive" formulation, along the lines of what a dynamic programming principle (DPP) ought to be. However, we show first that the causal quasi-Markov transport plans we will introduce, and their associated optimal transport problem, do possess a DPP. This fact (Theorem 4.2) together with Proposition 4.5 will then allow us to prove Theorem 2.6.

**Definition 4.1.** *A causal transport plan* $\gamma \in \Pi_c(\mu, \nu)$ *is called causal quasi-Markov, which we denote by* $\gamma \in \Pi_{cqm}(\mu, \nu)$, *if* $\gamma^{x_1, \ldots, x_t, y_1, \ldots, y_t}(dx_{t+1}, dy_{t+1}) = \gamma^{x_t, y_1, \ldots, y_t}(dx_{t+1}, dy_{t+1})$ *for every* $t$.

Of course $\Pi_{cqm}(\mu, \nu) \neq \emptyset$ iff $\mu$ is a Markov measure, in which case $\gamma^{x_1, \ldots, x_t, y_1, \ldots, y_t}(dx_{t+1}) = \mu^{x_t}(dx_{t+1})$ must hold. We introduce the causal quasi-Markov transport problem:

(Pcqm)
$$\inf_{\gamma \in \Pi_{cqm}(\mu,\nu)} \int c d\gamma$$

We now prove the DPP for (Pcqm):

**Theorem 4.2** (**DPP for** (Pcqm))**.** *Let* $\mu, \nu \in \mathcal{P}(\mathbb{R}^N)$, *suppose that* $\mu$ *is Markov and that the cost is semiseparable in the sense that* $c = \sum_t c_t(x_t, y_1, \ldots, y_t)$ *for non-negative Borel functions* $c_t$. *Then starting from* $V_N^c := 0$ *and defining recursively for* $t = N, \ldots, 2$:

(4.2) $\quad V_{t-1}^c(y_1, \ldots, y_{t-1}; m(dx_{t-1})) :=$
$$\inf_{\gamma \in \Pi\left(\int_{x_{t-1}} m(dx_{t-1})\mu^{x_{t-1}}(dx_t), \, \nu^{y_1, \ldots, y_{t-1}}(dy_t)\right)} \int \gamma(dx_t, dy_t) \Big\{ c_t(x_t, y_1, \ldots, y_t)$$
$$+ V_t^c(y_1, \ldots, y_t; \gamma^{y_t}(dx_t)) \Big\},$$

*we have that*

$$V_0^c := \inf_{\gamma \in \Pi(p_*^1 \mu, p_*^1 \nu)} \int \gamma(dx_1, dy_1) \{ c_1(x_1, y_1) + V_1^c(y_1; \gamma^{y_1}(dx_1)) \} = value(\text{Pcqm}),$$

*where each function* $V_{t-1}^c$ *is jointly universally measurable and convex in its last component. If moreover, Assumption 2.4 holds and each* $c_t$ *is l.s.c then* $V_{t-1}^c(y_1, \ldots, y_{t-1}; \cdot)$ *is l.s.c. and the problem* (Pcqm) *is attained.*

To be precise, we shall prove that $V_{t-1}^c$ is lower semianalytic (i.e. $\{V_{t-1}^c < r\}$ is analytic for each $r$; see [BS78, Definition 7.21]), which implies universal measurability.

*Proof.* We split the proof in several steps. It is clear from its definition that $V_t^c$ (for any $t$) is convex in its last component and generally bounded from below.

*STEP 1:* We first show that the sets:

$$D_{t-1} := \left\{ (y_1, \ldots, y_{t-1}, m, \gamma) \text{ s.t. } \gamma \in \Pi\left( \int_{x_{t-1}} m(dx_{t-1})\mu^{x_{t-1}}(dx_t), \, \nu^{y_1,\ldots,y_{t-1}}(dy_t) \right) \right\},$$

are Borel, so in particular analytic. Indeed, observe first that

$$\bar{D}_{t-1} := \left\{ (p, q, \gamma) \text{ s.t. } \gamma \in \Pi\left( q(dx_t), \, p(dy_t) \right) \right\},$$

is closed w.r.t. weak convergence of probability measures. So if we denote by $\Psi$ the map

$$(y_1, \ldots, y_{t-1}, m, \gamma) \mapsto \left( \nu^{y_1,\ldots,y_{t-1}}(dy_t), \int_{x_{t-1}} m(dx_{t-1})\mu^{x_{t-1}}(dx_t), \gamma \right),$$

we have $D_{t-1} = \Psi^{-1}(\bar{D}_{t-1})$ and so it suffices to show that the map $m(dx_{t-1}) \mapsto \int_{x_{t-1}} m(dx_{t-1})\mu^{x_{t-1}}(dx_t)$ is Borel, which further boils down to proving that if $g_1, \ldots, g_r$ are bounded Borel functions, then the sets

(4.3) $$\left\{ m \text{ Borel prob. measures s.t. } \int g_1 dm \geq 0, \ldots, \int g_r dm \geq 0 \right\},$$

are Borel. But this is now a straightforward monotone class argument which we omit.

For future use in Step 5, we observe that the sets

$$\left\{ (m, \gamma) \text{ s.t. } \gamma \in \Pi\left( \int_{x_{t-1}} m(dx_{t-1})\mu^{x_{t-1}}(dx_t), \, \nu^{y_1,\ldots,y_{t-1}}(dy_t) \right) \right\},$$

which are generally Borel only (as a fibers of the set $D_{t-1}$), are further closed as soon as Assumption 2.4 holds.

*STEP 2:* Let us now prove that the recursion is well-defined; namely, that the involved integrated functions are defined at most except for a null-set w.r.t. the integral. To be precise, we will show that these functions are lower semianalytic, and so universally measurable, thus in particular their integrals are well defined w.r.t. any Borel measure. For that matter, we shall first prove that

$$y_1, \ldots, y_t, m(dx_t) \mapsto V_t^c(y_1, \ldots, y_t; m)$$

is lower semianalytic (l.s.a. in short). By [Kal02, Lemma 1.40], we know that for the regular kernel of a given $\gamma(dx_t, dy_t)$, the application $y_t \mapsto \gamma^{y_t}(dx_t)$ is Borel measurable, then so is the map $(y_1, \ldots, y_t, \gamma) \mapsto (y_1, \ldots, y_t, \gamma^{y_t}(dx_t))$. Therefore, by [BS78, Lemma 7.30(3)] (on the composition of l.s.a. and Borel maps) we deduce that $(y_1, \ldots, y_t, \gamma) \mapsto V_t^c(y_1, \ldots, y_t; \gamma^{y_t})$ is l.s.a for $\gamma$ participating in the infimum in (4.2) and in particular the integral is indeed well-defined.

We start with $V_{N-1}^c$. The cost

$$(y_1, \ldots, y_{N-1}, m, \gamma) \mapsto \int \gamma(dx_N, dy_N)c_N(x_N, y_1, \ldots, y_N),$$

is Borel measurable and so in particular l.s.a. To see this, take $g_1 = -[c_N \wedge s]$ in (4.3) and take $s \to +\infty$. We can write $V_{N-1}^c(y_1, \ldots, y_{N-1}; m)$ as the infimum of this cost over the fiber of the set $D_{N-1}$ at $(y_1, \ldots, y_{N-1}, m)$. By [BS78, Proposition 7.47] and Step 1, the function $V_{N-1}^c$ is l.s.a. on the projection of $D_{N-1}$ onto its $(y_1, \ldots, y_{N-1}, m)$ components. By reverse induction, let us suppose that $V_t^c$ is l.s.a. and prove that $V_{t-1}^c$ is likewise. The cost to consider is now

$$(y_1, \ldots, y_{t-1}, m, \gamma) \mapsto \int \gamma(dx_t, dy_t)c_t(x_t, y_1, \ldots, y_t) + \int \nu^{y_1,\ldots,y_{t-1}}(dy_t)V_t^c(y_1, \ldots, y_t; \gamma^{y_t}),$$

the first term of which is Borel as before and whose second term is l.s.a. by virtue of the inductive step, the discussion about composition of l.s.a. and Borel maps above, and by [BS78, Proposition 7.48] on the integration of l.s.a. functions with respect to Borel kernels. By [BS78, Lemma 7.30(4)] we get that their sum is l.s.a. and so by [BS78, Proposition 7.47] and Step 1 again we conclude that $V_{t-1}^c$ is l.s.a.

*STEP 3:* By the previous point, and the way we wrote (4.2) as an infimum of a l.s.a. cost over a fiber of an analytic set in Step 2, we can perform for any $\varepsilon > 0$ by [BS78, Proposition 7.50(b)] a selection

$$(y_1, \ldots, y_{t-1}, m) \mapsto L_{t-1,\varepsilon}^{y_1,\ldots,y_{t-1};m}(dx_t, dy_t) \in \Pi\left( \int_{x_{t-1}} m(dx_{t-1})\mu^{x_{t-1}}(dx_t), \, \nu^{y_1,\ldots,y_{t-1}}(dy_t) \right),$$

so that the above mapping is universally measurable (i.e. measurable w.r.t. any completion of the corresponding Borel sets of its domain) and

$$V_{t-1}^c(y_1, \ldots, y_{t-1}, m) + \varepsilon \geq \int L_{t-1,\varepsilon}^{y_1,\ldots,y_{t-1};m}(dx_t, dy_t)\big[c_t(x_t, y_1, \ldots, y_t)$$
$$+ V_t^c(y_1, \ldots, y_t; (L_{t-1,\varepsilon}^{y_1,\ldots,y_{t-1};m})^{y_t})\big],$$

so each $L$ is an $\varepsilon$-optimizer for the corresponding problem. We will now build a measure, whose successive kernels will solve the recursions (4.2) at each step, modulo an $\varepsilon$ margin. Start with any $\varepsilon$-optimizer $\gamma_{0,\varepsilon}(dx_1, dy_1)$ of $V_0^c$. Then take $y_1 \mapsto \gamma_{1,\varepsilon}^{y_1}(dx_2, dy_2) := L_{1,\varepsilon}^{y_1;\gamma_0^{y_1}}(dx_2, dy_2)$, which is universally measurable by the composition result [BS78, Proposition 7.44]. Inductively, if $(y_1, \ldots, y_{t-1}) \mapsto \gamma_{t-1,\varepsilon}^{y_1,\ldots,y_{t-1}}(dx_t, dy_t)$ has been constructed in a universally measurable way, we build

$$(y_1, \ldots, y_t) \mapsto \gamma_{t,\varepsilon}^{y_1,\ldots,y_t}(dx_{t+1}, dy_{t+1}) := L_{t,\varepsilon}^{y_1,\ldots,y_t;(\gamma_{t-1,\varepsilon}^{y_1,\ldots,y_{t-1}})^{y_t}}(dx_{t+1}, dy_{t+1}),$$

which is universally measurable by [BS78, Proposition 7.44], the inductive step, and the fact that the function that associates a regular kernel to a product measure is Borel. This finishes the induction, and by construction

$$(4.4) \qquad \gamma_{t,\varepsilon}^{y_1,\ldots,y_t}(dx_{t+1}, dy_{t+1}) \in \Pi\left(\int_{x_t}(\gamma_{t-1,\varepsilon}^{y_1,\ldots,y_{t-1}})^{y_t}(dx_t)\mu^{x_t}(dx_{t+1}), \nu^{y_1,\ldots,y_t}(dy_{t+1})\right),$$

and $\gamma_{t,\varepsilon}^{y_1,\ldots,y_t}$ attains $V_t^c(y_1, \ldots, y_t; (\gamma_{t-1,\varepsilon}^{y_1,\ldots,y_{t-1}})^{y_t})$ except for an $\varepsilon$ margin. We now define

$$(x_t, y_1, \ldots, y_t) \mapsto \Gamma_{t,\varepsilon}^{x_t,y_1,\ldots,y_t}(dx_{t+1}, dy_{t+1}) := \mu^{x_t}(dx_{t+1})(\gamma_{t,\varepsilon}^{y_1,\ldots,y_t})^{x_{t+1}}(dy_{t+1}),$$

which is universally measurable by [BS78, Proposition 7.44] again. By [BS78, Proposition 7.45] the successive concatenation of these kernel induce a unique Borel measure $\Gamma_\varepsilon$ such that

$$\Gamma_\varepsilon(dx_1, \ldots, dx_N, dy_1, \ldots, dy_N) = \gamma_{0,\varepsilon}(dx_1, dy_1)\Gamma_{1,\varepsilon}^{x_1,y_1}(dx_2, dy_2)\ldots$$
$$\Gamma_{t,\varepsilon}^{x_t,y_1,\ldots,y_t}(dx_{t+1}, dy_{t+1})\ldots\Gamma_{N-1,\varepsilon}^{x_{N-1},y_1,\ldots,y_{N-1}}(dx_N, dy_N).$$

By construction, this measure is causal quasi-Markov and its $x$-marginal is exactly $\mu$. As for the $y$-marginal, it is easy to see that $\Gamma_\varepsilon(dy_1) = \nu(y_1)$ and that

$$\Gamma_\varepsilon(dy_1, dy_2) = \int_{x_1,x_2} \gamma_{0,\varepsilon}(dx_1, dy_1)\mu^{x_1}(dx_2)(\gamma_{1,\varepsilon}^{y_1})^{x_2}(dy_2)$$
$$= \nu(dy_1)\int_{x_2}\left(\int_{x_1}\gamma_{0,\varepsilon}(dx_1)\mu^{x_1}(dx_2)\right)(\gamma_{1,\varepsilon}^{y_1})^{x_2}(dy_2) = \nu(dy_1)\nu^{y_1}(dy_2),$$

observing (from (4.4)) in the last line, inside the brackets, that this is exactly the first marginal of $\gamma_{1,\varepsilon}^{y_1}$. Inductively, one can verify that $\Gamma_\varepsilon$ has $\nu$ as its $y$-marginal.

*STEP 4:* We now prove the equality $V_0^c = value(\mathrm{Pcqm})$. By the previous step, $V_0^c + N\varepsilon \geq value(\mathrm{Pcqm})$, since $\Gamma_\varepsilon$ is feasible for (Pcqm) and because by construction $\Gamma_\varepsilon$ was designed $\varepsilon$-optimal at each step, so that overall $V_0^c + N\varepsilon \geq \int c\, d\Gamma_\varepsilon$. By letting $\varepsilon \to 0$, we get $V_0^c \geq value(\mathrm{Pcqm})$. On the other hand, given any $\gamma \in \Pi_{cqm}(\mu, \nu)$, we clearly have that

$$(4.5) \qquad \gamma^{y_1,\ldots,y_t}(dx_{t+1}, dy_{t+1}) \in \Pi\big(\int_{x_t}(\gamma^{y_1,\ldots,y_{t-1}})^{y_t}(dx_t)\mu^{x_t}(dx_{t+1}),\, \nu^{y_1,\ldots,y_t}(dy_{t+1})\big),$$

and further we can write

$$\int c_t\, d\gamma = \int \gamma(dy_1, \ldots, dy_{t-1})\gamma^{y_1,\ldots,y_{t-1}}(dx_t, dy_t)c_t(x_t, y_1, \ldots, y_t) =$$
$$\int \gamma(dx_1, dy_1)\gamma^{y_1}(dx_2, dy_2)\gamma^{y_1,y_2}(dx_3, dy_3)\ldots\gamma^{y_1,\ldots,y_{t-1}}(dx_t, dy_t)c_t(x_t, y_1, \ldots, y_t),$$

so

$$(4.6) \qquad \int c\, d\gamma = \int \gamma(dx_1, dy_1)\left\{c_1 + \int \gamma^{y_1}(dx_2, dy_2)\left\{c_2 + \int \gamma^{y_1,y_2}(dx_3, dy_3)\{c_3 + \ldots\}\right\}\right\}.$$

Combining this with (4.5) we get $value(\mathrm{Pcqm}) \geq V_0^c$ and conclude the desired equality.

From now on we take Assumption 2.4 for granted and assume l.s.c. costs.

*STEP 5:* We establish now that $V_{t-1}^c$ is lower semicontinuous in its last argument (i.e. $m$), the other ones being fixed. The desired lower semicontinuity could be proved by hand using convex combination of kernels as in Step 6 below, but a simpler argument is to invoke a "maximum theorem" as in [BS78,

Proposition 7.33], which establishes the lower semicontinuity of a value function as long as the cost is jointly lower semicontinuous (in our case w.r.t. $(m, \gamma)$, since the $y$'s are fixed, and this clearly holds) and as soon as the constraint set is a fiber of a closed set (which is true by the last observation in Step 1). The lower semicontinuity in $(m, \gamma)$ is proved by reverse induction, much as in Step 2, and for [BS78, Proposition 7.33] one can assume that the $\gamma$'s live a priori in a compact space, since varying the $m$'s in a tight set only will "move" the constraint set $\Pi$ inside a larger tight set in the product.

*STEP 6:* We now prove that each of the problems in (4.2) is attained. The $\int c_t d\gamma$ part of the cost is l.s.c. and linear, so it causes no trouble and we may assume $c_t \equiv 0$. We take a minimizing sequence $\gamma_n$ for (4.2):

$$V_{t-1}^c(y_1, \ldots, y_{t-1}; m) = \lim_n \int \nu^{y_1, \ldots, y_{t-1}}(dy_t) V^c(y_1, \ldots, y_t; \gamma_n^{y_t}).$$

Trivially the sequence $\int_{y_t} \nu^{y_1, \ldots, y_{t-1}}(dy_t) \gamma_n^{y_t}(dx_t)$ is tight, since it is identically equal to the measure $\int_{x_{t-1}} m(dx_{t-1}) \mu^{x_{t-1}}(dx_t)$. By [Bal, Theorem 3.15] we conclude that there is a subsequence of the kernel $\{y_t \mapsto \gamma_n^{y_t}\}_n$, which we denote the same, and a kernel $\{y_t \mapsto \gamma^{y_t}\}$ such that the sequence of Césaro averages

$$y_t \mapsto \tilde{\gamma}_n^{y_t} := \frac{1}{n} \sum_{i \leq n} \gamma_i^{y_t}$$

converges weakly, for $\nu^{y_1, \ldots, y_{t-1}}(dy_t)$-a.e. $y_t$, to the kernel $y_t \mapsto \gamma^{y_t}$ (i.e. as measures on $x_t$, see [Bal, Definition 3.10]) and further it holds that $\int \nu^{y_1, \ldots, y_{t-1}}(dy_t) \gamma^{y_t}(dx_t)$ equals $\int_{x_{t-1}} m(dx_{t-1}) \mu^{x_{t-1}}(dx_t)$ by [Bal, Corollary 3.14], so the measure $\gamma := \nu^{y_1, \ldots, y_{t-1}}(dy_t) \gamma^{y_t}(dx_t)$ is feasible for (4.2). All in all it follows

$$V_{t-1}^c(y_1, \ldots, y_{t-1}; m) = \liminf_n \frac{1}{n} \sum_{i \leq n} \int \nu^{y_1, \ldots, y_{t-1}}(dy_t) V^c(y_1, \ldots, y_t; \gamma_i^{y_t})$$

$$\geq \liminf_n \int \nu^{y_1, \ldots, y_{t-1}}(dy_t) V^c(y_1, \ldots, y_t; \tilde{\gamma}_n^{y_t})$$

$$\geq \int \nu^{y_1, \ldots, y_{t-1}}(dy_t) \liminf_n V^c(y_1, \ldots, y_t; \tilde{\gamma}_n^{y_t})$$

$$= \int \nu^{y_1, \ldots, y_{t-1}}(dy_t) V^c(y_1, \ldots, y_t; \gamma^{y_t}),$$

by convexity, Fatou's Lemma (remember the $V^c$'s are bounded below) and the lower semicontinuity established in the previous Step. Therefore, $\gamma$ is feasible and optimal.

*Step 7:* By the previous step, we may now go back to Step 3 and find by [BS78, Proposition 7.50(b)] a selection of optimizers at each time (as opposed to only $\varepsilon$-optimizers) $L_t^{y_1, \ldots, y_t; m}$. Then we can redo Step 3 with $\varepsilon = 0$, building a global measure $\Gamma$ which exactly solves the recursion and is feasible for (Pcqm), and so is optimal for it. $\qquad\square$

Some observations regarding the previous theorem and its proof:

**Remark 4.3.** *Even if c does not have the stated semiseparable structure, a look at the proof shows that the recursion is well-posed and that its value gives an upper bound for value(Pcqm). The semiseparable structure is crucial for the lower bound (see (4.6)) only.*

**Remark 4.4.** *As a by-product of the previous proof, we see from Step 6 therein a way to prove attainability of general transport problems as in [GRST15] in the presence of lower semicontinuity and convexity of the cost w.r.t. the kernel, without the assumption that the state space be compact or discrete-like.*

As we have mentioned in Theorem 2.6, each optimization problem in DPP is of the form of general (non-linear) transports on the line. Under Assumption 3.2 and for the lower semicontinuous cost function each of the value functions in (4.2) is convex in the kernel and jointly l.s.c. Assuming additionally that all the regular kernels of measures $\mu, \nu$ are compactly supported, we can use the duality result of Gozlan et al. [GRST15, Theorem 9.5]. Denote $\eta(dx_t) = \int_{x_{t-1}} m(dx_{t-1}) \mu^{x_{t-1}}(dx_t)$, and

$$\bar{c}_t^{y_1, \ldots, y_{t-1}}(y_t; p) = \int_{x_t} c_t(x_t, y_1, \ldots, y_t) p(dx_t) + V_t^c(y_1, \ldots, y_{t-1}, y_t; p(dx_t)).$$

Then the following dual formulation holds:

$$(4.7) \qquad V_{t-1}^c(y_1, \ldots, y_{t-1}; m(dx_{t-1})) = \sup_\phi \left\{ \int R_{\bar{c}_t}^{y_1, \ldots, y_{t-1}} \phi(y_t) \nu^{y_1, \ldots, y_{t-1}}(dy_t) - \int \phi(x_t) \eta(dx_t) \right\},$$

where

$$R_{\bar{c}_t^{y_1,\ldots,y_{t-1}}} \phi(y_t) = \inf_{p \in \mathcal{P}(\mathbb{R})} \left\{ \int \phi(x_t)p(dx_t) + \bar{c}_t^{y_1,\ldots,y_{t-1}}(y_t;p) \right\}.$$

We finally establish a sufficient condition for (Pcqm) and (Pc) to be equivalent:

**Proposition 4.5.** *Let $\mu$ be Markov and the cost be semiseparable. Then $value(\text{Pcqm}) = value(\text{Pc})$ and if (Pc) is attained, then there is some optimal causal transport which is further causal quasi-Markov.*

*Proof.* We prove that under the given assumptions, any causal plan has a related causal quasi-Markov plan which incurs in the same cost. Let $\gamma \in \Pi_c(\mu,\nu)$, and build

$$d\hat{\gamma} = \gamma(dx_1,dy_1)\gamma^{x_1,y_1}(dx_2,dy_2)\gamma^{x_2,y_1,y_2}(dx_3,dy_3)\ldots\gamma^{x_{N-1},y_1,\ldots,y_{N-1}}(dx_N,dy_N).$$

As $\hat{\gamma}^{x_1,\ldots,x_t,y_1,\ldots,y_t} = \gamma^{x_t,y_1,\ldots,y_t} = \hat{\gamma}^{x_t,y_1,\ldots,y_t}$ we have $\hat{\gamma}$ is causal and its $x$-marginal is $\mu$. Indeed, $\forall t$:

$$\hat{\gamma}^{x_1,\ldots,x_t,y_1,\ldots,y_t}(dx_{t+1}) = \gamma^{x_t,y_1,\ldots,y_t}(dx_{t+1})$$
$$= \int_{x_1,\ldots,x_{t-1}} \gamma^{x_1,\ldots,x_t,y_1,\ldots,y_t}(dx_{t+1})\gamma^{x_t,y_1,\ldots,y_t}(dx_1,\ldots,dx_{t-1})$$
$$= \int_{x_1,\ldots,x_{t-1}} \mu^{x_t}(dx_{t+1})\gamma^{x_t,y_1,\ldots,y_t}(dx_1,\ldots,dx_{t-1}) = \mu^{x_t}(dx_{t+1}).$$

We finally prove that for every function of the form $H := H(x_t,y_1,\ldots,y_t)$, holds that $\int H d\gamma = \int H d\hat{\gamma}$. This shows first that $\gamma$ and $\hat{\gamma}$ incur in the same cost under the semiseparability assumption, and second that the $y$-marginal of $\hat{\gamma}$ is $\nu$, which establishes $\hat{\gamma} \in \Pi_{cqm}(\mu,\nu)$.

$$\int H d\hat{\gamma} = \int H\gamma(dx_1,dx_2,dy_1,dy_2)\gamma^{x_2,y_1,y_2}(dx_3,dy_3)\hat{\gamma}^{x_3,y_1,y_2,y_3}(dx_4,\ldots,dy_4,\ldots)$$
$$= \int H\gamma(dx_2,dy_1,dy_2)\gamma^{x_2,y_1,y_2}(dx_3,dy_3)\hat{\gamma}^{x_3,y_1,y_2,y_3}(dx_4,\ldots,dy_4,\ldots)$$
$$= \int H\gamma(dx_2,dx_3,dy_1,dy_2,dy_3)\gamma^{x_3,y_1,y_2,y_3}(dx_4,dy_4)\hat{\gamma}^{x_4,y_1,y_2,y_3,y_4}(dx_5,\ldots,dy_5,\ldots)$$
$$= \int H\gamma(dx_3,dy_1,dy_2,dy_3)\gamma^{x_3,y_1,y_2,y_3}(dx_4,dy_4)\hat{\gamma}^{x_4,y_1,y_2,y_3,y_4}(dx_5,\ldots,dy_5,\ldots)$$
$$\ldots = \int H\gamma(dx_{t-1},dy_1,\ldots,dy_{t-1})\gamma^{x_{t-1},y_1,\ldots,y_{t-1}}(dx_t,dy_t) = \int H d\gamma.$$

$\square$

**Remark 4.6.** *With Proposition 4.5 we can conclude that, whenever $\mu$ is Markov and the cost is l.s.c. and semiseparable, the problem (Pcqm) has a solution. Indeed, since (Pc)=(Pcqm) and by [Las15, Corollary 1] the former is attained, one constructs also an optimizer for the latter. Our Theorem 4.2 yields the same with a stronger assumption in a self-contained way, through DPP. Let us stress that the causal quasi-Markov constraint is not a linear/convex one.*

All in all, the proof of Theorem 2.6 is now trivial:

*Proof of Theorem 2.6.* By Proposition 4.5 $value(\text{Pcqm}) = value(\text{Pc})$ and by Theorem 4.2 we have the recursion and all the other properties. $\square$

We close the present part with a preliminary illustration of how the causal DPP can be used in practice; here we show a condition giving the equivalence between (Pc) and (Pbc):

**Corollary 4.7.** *Consider $N = 2$ and a separable cost of the form $c = c_1(x_1,y_1) + |x_2 - y_2|$, with $c_1$ l.s.c. and bounded from below. Suppose that for every $z,y_1$ one of the following two cases holds*

$$\text{either } [\forall x_1 : F_{\mu^{x_1}}(z) \geq F_{\nu^{y_1}}(z)] \quad \text{or} \quad [\forall x_1 : F_{\mu^{x_1}}(z) \leq F_{\nu^{y_1}}(z)],$$

*Then (Pc) is equivalent to (Pbc).*

*Proof.* Borrowing from Theorem 2.6, call

$$V_1(y_1,\gamma^{y_1}) = \inf_{m \in \Pi(\int \gamma^{y_1}(dx_1)\mu^{x_1},\nu^{y_1})} \int |x_2 - y_2|m(dx_2,dy_2),$$

which by e.g. [Vil03, eq.(2.48), p.75] equals $\int_z \left| F_{\int \gamma^{y_1}(dx_1)\mu^{x_1}}(z) - F_{\nu^{y_1}}(z) \right| dz$. But $F_{\int \gamma^{y_1}(dx_1)\mu^{x_1}}(z) = \int \gamma^{y_1}(dx_1)F_{\mu^{x_1}}(z)$, so by our technical assumption:

$$\int_z \left| F_{\int \gamma^{y_1}(dx_1)\mu^{x_1}}(z) - F_{\nu^{y_1}}(z) \right| dz = \int_z \left| \int_{x_1} \gamma^{y_1}(dx_1)F_{\mu^{x_1}}(z) - F_{\nu^{y_1}}(z) \right| dz$$
(4.8)
$$= \int_z \int_{x_1} \gamma^{y_1}(dx_1) \left| F_{\mu^{x_1}}(z) - F_{\nu^{y_1}}(z) \right| dz,$$

and by Fubini's Theorem we get $V_1(y_1, \gamma^{y_1}) = \int_{x_1} \gamma^{y_1}(dx_1) \inf_{m \in \Pi(\mu^{x_1}, \nu^{y_1})} \int |x_2 - y_2| dm$ and again from Theorem 2.6 we obtain that

$$Value(\text{Pc}) = \inf_{\gamma \in \Pi(p^1_* \mu, p^1_* \nu)} \int \gamma(dx_1, dy_1) \left\{ c_1(x_1, y_1) + \inf_{m \in \Pi(\mu^{x_1}, \nu^{y_1})} \int |x_2 - y_2| dm \right\},$$

which as we shall see in Proposition 5.2, equals the bicausal DPP, so we conclude. $\qquad\square$

## 5. The bicausal case

We can obtain a similar DPP for bicausal transport plans. Let us introduce:

(Dyn-Pbc) $\quad \inf_{\gamma^1 \in \Pi(p^1_* \mu, p^1_* \nu)} \int \gamma^1(dx_1, dy_1) \inf_{\gamma^2 \in \Pi(\mu^{x_1}, \nu^{y_1})} \int \gamma^2(dx_2, dy_2) \dots$

$$\dots \inf_{\gamma^N \in \Pi(\mu^{x_1, \dots, x_{N-1}}, \nu^{y_1, \dots, y_{N-1}})} \int \gamma^N(dx_N, dy_N) c(x_1, \dots, x_N, y_1, \dots, y_N).$$

The previous recursive problem is motivated by the following structure result, the proof of which is analogous to the causal case, so we omit it:

**Proposition 5.1.** *Let $\mu, \nu \in \mathcal{P}(\mathbb{R}^N)$.*
*If $\gamma \in \Pi_{bc}(\mu, \nu)$ is decomposed as in (2.1), then the following conditions on the kernels hold:*
  (i) *$\bar{\gamma} \in \Pi(p^1_* \mu, p^1_* \nu)$, and*
  (ii) *successively for $t < N$ and for $\gamma$-almost every $x_1, \dots, x_t, y_1, \dots, y_t$ holds*

$$\gamma^{x_1, \dots, x_t, y_1, \dots, y_t}(dx_{t+1}, dy_{t+1}) \in \Pi(\mu^{x_1, \dots, x_t}(dx_{t+1}), \nu^{y_1, \dots, y_t}(dy_{t+1})).$$

*Conversely, given regular kernels*

$$\bar{\gamma}(dx_1, dy_1), \gamma^{x_1, y_1}(dx_2, dy_2), \dots, \gamma^{x_1, \dots, x_{N-1}, y_1, \dots, y_{N-1}}(dx_N, dy_N),$$

*satisfying the properties $(i) - (ii)$, the measure $\gamma$ constructed as in (2.1) belongs to $\Pi_{bc}(\mu, \nu)$.*

The recursion corresponding to (Dyn-Pbc) (starting from $V^c_N := c$) is:

(5.1) $\quad V^c_t(x_1, \dots, x_t, y_1, \dots, y_t) =$

$$\inf_{\gamma^{t+1} \in \Pi(\mu^{x_1, \dots, x_t}, \nu^{y_1, \dots, y_t})} \int \gamma^{t+1}(dx_{t+1}, dy_{t+1}) V^c_{t+1}(x_1, \dots, x_{t+1}, y_1, \dots, y_{t+1}),$$

and so we want to compare the values of (Dyn-Pbc), (Pbc) and

$$V^c_0 := \inf_{\gamma^1 \in \Pi(p^1_* \mu, p^1_* \nu)} \int V^c_1(x_1, y_1) \gamma^1(dx_1, dy_1).$$

We now give the DPP for the bicausal case, which appeared in [Rüs85, Theorem 3]; we will prove it with our methods, obtaining further attainability and some regularity of the value function. Observe that no Markovianity of $\mu$ or separability of the costs is needed, and that the above value function (and recursion) are more tractable than in the causal quasi-Markov case (4.2).

**Proposition 5.2.** *Given a Borel bounded from below cost function $c$, we have that the recursive optimization problem (Dyn-Pbc) is well-defined, namely the successive integrals in (5.1) are well-defined, and the values of (Dyn-Pbc), (Pbc) and $V^c_0$ coincide. If further $c$ is l.s.c. and Assumption 3.2 holds, then there is a bicausal optimizer and the value functions $V^c_t$ are all l.s.c.*

*Proof.* Since $(x_1, \dots, x_t, y_1, \dots, y_t) \mapsto (\mu^{x_1, \dots, x_t}, \nu^{y_1, \dots, y_t})$ is Borel and the set $\{(p, q, \gamma) : \gamma \in \Pi(p, q)\}$ is clearly closed, we get that

$$D_t := \{(x_1, \dots, x_t, y_1, \dots, y_t, \gamma) : \gamma \in \Pi(\mu^{x_1, \dots, x_t}, \nu^{y_1, \dots, y_t})\},$$

is Borel. As in the proof of Theorem (4.2), we start by proving recursively that the $V^c_t$'s are lower semianalytic (l.s.a). As in Step 2 therein one observes first that

$$(x_1, \dots, x_{N-1}, y_1, \dots, y_{N-1}, \gamma) \mapsto \int \gamma(dx_N, dy_N) c(x_1, \dots, x_N, y_1, \dots, y_N),$$

is Borel, and so l.s.a. Since $V^c_{N-1}(x_1, \dots, x_{N-1}, y_1, \dots, y_{N-1})$ is the infimum of this function over the fiber of $D_{N-1}$ at $(x_1, \dots, x_{N-1}, y_1, \dots, y_{N-1})$ we get by [BS78, Proposition 7.47] that $V^c_{N-1}$ is l.s.a. and in particular universally measurable. The recursive step is obvious, and we get this result for each $V^c_t$, and so the integrals in (5.1) are well-defined. Now take, as in Step 3 of the proof of Theorem (4.2), a universally measurable selection of $\varepsilon$-optimizers for $V^c_t$; call these $\gamma^{x_1, \dots, x_t, y_1, \dots, y_t}_{t, \varepsilon}$. Build then $\Gamma_\varepsilon$ as the unique Borel measure that comes out of concatenating these ([BS78, Proposition 7.45]). By construction

each $\Gamma_\varepsilon$ is in $\Pi(\mu,\nu)$ (see Proposition 5.1) and $V_0^c + N\varepsilon \geq \int c d\Gamma_\varepsilon \geq value(\text{Pbc})$ and so $V_0^c \geq value(\text{Pbc})$. The reverse inequality follows trivially from Proposition 5.1.

If Assumption 3.2 holds then $D_t$ is closed and so if further $c$ is l.s.c. we argue as in Step 5 of the proof of Theorem (4.2), proving that the value functions are l.s.c. as well. In turn, applying now [BS78, Proposition 7.50(b)] we can obtain a universally measurable selection of optimizers $\gamma_t^{x_1,\ldots,x_t,y_1,\ldots,y_t}$ and with them we build an optimal bicausal transport $\Gamma$ (just like the previous paragraph with $\varepsilon = 0$). □

Since at each step of the dynamic programming principle (Dyn-Pbc), or equivalently (5.1), we have a usual optimal transport problem, we can write these in their dual formulation; for simplicity we assume that $c$ is l.s.c. and non-negative (see however [BS11] for an extension). In effect, the value function at time $t$ for each $t = N-1, N-2, \ldots, 1$ is obtained recursively as:

$$V^c(t; x_1, \ldots, x_t, y_1 \ldots, y_t) :=$$

$$\sup_{\substack{\phi_{t+1}, \psi_{t+1} \in C_b(\mathbb{R}),\\ \phi_{t+1}(x_{t+1}) + \psi_{t+1}(y_{t+1}) \leq V^c(t+1; x_1, \ldots, x_{t+1}, y_1, \ldots, y_{t+1})}} \left\{ \int \phi_{t+1}(x_{t+1}) \mu^{x_1,\ldots,x_t}(dx_{t+1}) \right.$$

$$\left. + \int \psi_{t+1}(y_{t+1}) \nu^{y_1,\ldots,y_t}(dy_{t+1}) \right\},$$

so the value of (Dyn-Pbc) is also given by

$$V^c(0) := \sup_{\substack{\phi_1, \psi_1 \in C_b(\mathbb{R})\\ \phi_1(x_1) + \psi_1(y_1) \leq V^c(1; x_1, y_1)}} \int \phi_1(x_1) p_*^1 \mu(dx_1) + \int \psi_1(y_1) p_*^1 \nu(dy_1).$$

Observe that this recursive structure of the dual problem is not obvious from (Dbc), but can be guessed a posteriori thanks to the apparent "primal" recursive structure. A simpler picture of (Dyn-Pbc) arises when $c$ has separable structure as in (2.10); see [PP12, Theorem 7.2].

We give now the belated proof of Theorem 2.7:

*Proof of Theorem 2.7.* The result follows from Proposition 5.3 below, if the causal/bicausal equality is proved, which we do now. Start with $\gamma \in \Pi_c(\mu,\nu)$ and decompose it as in (2.1). From Proposition 2.3, we know that the following conditions ($t < N$) are satisfied by the kernels $\bar\gamma(dx_1, dy_1), \gamma^{x_1,y_1}(dx_2, dy_2)$, up to $\gamma^{x_1,\ldots,x_{N-1},y_1,\ldots,y_{N-1}}(dx_N, dy_N)$:

$$(5.2) \qquad \gamma^{x_1,\ldots,x_t,y_1,\ldots,y_t}(dx_{t+1}, \mathbb{R}) = \mu_{t+1}(dx_{t+1})$$

and

$$(5.3) \qquad \int_{x_1,\ldots,x_t} \gamma^{x_1,\ldots,x_t,y_1,\ldots,y_t}(\mathbb{R}, dy_{t+1}) \gamma^{y_1,\ldots,y_t}(dx_1, \ldots, dx_t) \; = \; \nu^{y_1,\ldots,y_t}(dy_{t+1})$$
$$= \gamma^{y_1,\ldots,y_t}(\mathbb{R}, dy_{t+1})$$

We can rewrite (5.2) in the following way:

$$(5.4) \qquad \int_{x_1,\ldots,x_t} \gamma^{x_1,\ldots,x_t,y_1,\ldots,y_t}(dx_{t+1}, \mathbb{R}) \gamma^{y_1,\ldots,y_t}(dx_1, \ldots, dx_t) \; = \; \mu_{t+1}(dx_{t+1})$$
$$= \gamma^{y_1,\ldots,y_t}(dx_{t+1}, \mathbb{R})$$

Therefore, we can construct a new plan $\tilde\gamma$ as follows

$$\tilde\gamma(dx_1, \ldots, dx_N, dy_1, \ldots, dy_N) = \bar\gamma(dx_1, dy_1) \gamma^{y_1}(dx_2, dy_2) \ldots \gamma^{y_1,\ldots,y_{N-1}}(dx_N, dy_N),$$

with

$$(5.5) \qquad \gamma^{y_1,\ldots,y_t}(dx_{t+1}, dy_{t+1}) = \int_{x_1,\ldots,x_t} \gamma^{x_1,\ldots,x_t,y_1,\ldots,y_t}(dx_{t+1}, dy_{t+1}) \gamma^{y_1,\ldots,y_t}(dx_1, \ldots, dx_t).$$

Due to (5.4) and (5.3), one can see that for any $t < N$ each kernel

$$\gamma^{y_1,\ldots,y_t} \in \Pi(\mu_{t+1}(dx_{t+1}), \nu^{y_1,\ldots,y_t}(dy_{t+1})).$$

Then, from Proposition (5.1), we know that $\tilde\gamma \in \Pi_{bc}(\mu,\nu)$. Writing down the minimization over the kernels (5.5), we have exactly (Pbc). Since the kernels of $\tilde\gamma$ and $\gamma$ are connected via (5.4) and (5.3), one can observe that the value of the bicausal transport problem is less or equal than the causal one, giving us the desired result. □

*Proof of Corollary 2.8.* Follows from Theorem 2.7 by observing that the map

$$(x_1, \ldots, x_N, y_1, \ldots, y_N) \mapsto (x_1, x_2 - x_1, \ldots, x_N - x_{N-1}, y_1, y_2 - y_1, \ldots, y_N - y_{N-1}),$$

preserves causality when applied to $\gamma \in \Pi_c(\mu, \nu)$. □

The proof of Theorem 2.7 rests on the following result, which needs Condition (5.6) encompassing at the same time the independence condition (2.11) and Rüschendorf's "monotone regression dependence" in [Rüs85, Corollary 2].

**Proposition 5.3.** *For each $t = 1, \ldots, N - 2$, all $(x_1, \ldots, x_t), (y_1, \ldots, y_t)$, and $u$, suppose*

$$(5.6) \qquad \left(F_{\mu^{x_1, \ldots, x_t, \bar{x}}}(u) - F_{\mu^{x_1, \ldots, x_t, x}}(u)\right)\left(G_{\nu^{y_1, \ldots, y_t, \bar{y}}}(u) - G_{\nu^{y_1, \ldots, y_t, y}}(u)\right) \geq 0,$$

*whenever $\bar{x} \geq x$ and $\bar{y} \geq y$. Assume further that*

$$c(x_1, \ldots, x_N, y_1, \ldots, y_N) := \sum_{t \leq N} c_t(x_t - y_t),$$

*where each $c_t$ is convex. Then a solution to (Pbc) is given by the Knothe-Rosenblatt rearrangement (2.8). Additionally, if $\mu$-a.s. all the conditional distributions of $\mu$ are atomless (e.g. if $\mu$ has a density), then this rearrangement is induced by the Monge map determined by (2.9).*

*Proof.* First, we give the proof for the case when the source measure is the product of its marginals, for $N = 2$. From Proposition 5.2 we know that the values of (Pbc) and (Dyn-Pbc) coincide, and by assumption $\mu(dx_1, dx_2) = \mu_1(dx_1)\mu_2(dx_2)$, so (Dyn-Pbc) has the form

$$(5.7) \quad \inf_{\gamma^1 \in \Pi(\mu_1, p_*^1\nu)} \int \gamma^1(dx_1, dy_1) \left[c_1(x_1 - y_1) + \inf_{\gamma^2 \in \Pi(\mu_2, \nu^{y_1})} \int c_2(x_2 - y_2)\gamma^2(dx_2, dy_2)\right].$$

From classical optimal transport (see [Vil03]) it is known that

$$\inf_{\gamma^2 \in \Pi(\mu_2, \nu^{y_1})} \int c_2(x_2 - y_2)\gamma^2(dx_2, dy_2) = \int c(F_{\mu_2}^{-1}(u_2) - F_{\nu^{y_1}}^{-1}(u_2))du_2.$$

The last value does not depend on $x_1$, therefore, it is constant for the minimization with respect to $\gamma^1$ in (5.7), and we conclude that the pair $(X_1^*, Y_1^*)$ will be the optimizer at this step, which in turn allow to construct $(X_2^*, Y_2^*)$ by the previous considerations. For general $N$ the result follows by induction, iterating the arguments so far. When $\mu$ is not the product of its marginals, yet the monotone regression assumption of [Rüs85] holds, the proof of this result is given in [Rüs85, Corollary 2]. Finally, it is easy to observe that Condition 5.6 unifies exactly the two cases described. □

We stress that in case the independent marginals condition 2.11 does not hold, and even if condition 5.6 is true, Example 7.2 shows that causal and bicausal values may differ.

The explicit form of the optimizers in Theorem 2.7 is obviously only true for $\mathbb{R}$-valued processes. Interestingly, if the transport problem consisted in mapping (in a bicausal way) $\mathbb{R}^L$-valued process, the recursive structure is just as in the case $L = 1$ and so if $\mu$ were the product of its marginals (each of them in $\mathcal{P}(\mathbb{R}^L)$ now) we would get a similar conclusion with the role of the monotone rearrangements taken by general Brenier's maps, under suitable conditions.

**Remark 5.4.** *We reassure the reader that the Knothe-Rosenblatt rearrangements in the form (2.8) are always bicausal, and in the Monge form (2.9) this is also the case as soon as all conditional distributions of the source measure are atomless. The argument for (2.8) is as follows: for all $g(\cdot, \cdot)$ there is a $G(\cdot)$ such that*

$$\int_0^1 \int_0^1 f(F_{\mu_1}^{-1}(u_1))g\left(F_{\nu_1}^{-1}(u_1), F_{\nu_{F_{\nu_1}^{-1}(u_1)}}^{-1}(u_2)\right)du_2 du_1 = \int_0^1 \int_0^1 f(F_{\mu_1}^{-1}(u_1))G \circ F_{\nu_1}^{-1}(u_1)du_1,$$

*so the law of $F_{\mu_1}^{-1}(U_1)$ given $F_{\nu_1}^{-1}(U_1)$ and $F_{\nu_{F_{\nu_1}^{-1}(U_1)}}^{-1}(U_2)$, equals the law of $F_{\mu_1}^{-1}(U_1)$ given $F_{\nu_1}^{-1}(U_1)$. The same holds inverting the roles of $\mu$ and $\nu$ and going to greater time indices. As for (2.9), the argument actually follows directly upon noticing for example that for $\mu$−a.e. $x_1$ the measure $[F_{\mu^{x_1}}]_*\mu^{x_1}$ is equal to Lebesgue measure on $[0, 1]$, under the given assumptions.*

5.1. **Digression into a classical functional inequality.** Another instance where Theorem 2.7 (and Condition 2.11) comes in handy, is the following interpretation of the proof of Talagrand's $\mathcal{T}_2$ inequality. First, recall (see [Tal96]): given a unit standard Gaussian measure $G^N$ and any other measure $\nu$ on $\mathbb{R}^N$, the following holds:

(T2) $$2Ent(\nu|G^N) \geq \mathcal{T}_2(G^N, \nu),$$

where $Ent$ denotes relative entropy and $\mathcal{T}_2(\cdot, \cdot)$ is the value of the optimal transport problem with quadratic cost. Equality holds iff $\nu$ is an affine translate of $G^N$. We establish here the related bicausal inequality

(CT2) $$2Ent(\nu|G^N) \geq \mathcal{T}_{bc,2}(G^N, \nu), \text{ for all } \nu \in \mathcal{P}(\mathbb{R}^N),$$

where $\mathcal{T}_{bc,2}(G^N, \nu)$ is the value of the optimal bicausal transport problem for quadratic cost. This clearly implies (T2) for $G^N$. It is the DPP for the bicausal transport problem that replaces the role of the usual tensorization trick; strictly speaking, the DPP is just giving a name to an intermediate step in Talagrand's well-known proof. The proof given here applies of course to other product measures, and in that setting is reminiscent of [Vil08, Proposition 22.5]; we stick to the gaussian case only for the sake of concreteness and because this relates to the Wiener case in continuous time.

**Proposition 5.5.** *For $G^N$ the unit standard Gaussian measure on $\mathbb{R}^N$, the bicausal transport-entropy inequality* (CT2) *holds.*

*Proof.* We start assuming (CT2) for $N = 1$, which holds by [Tal96], and prove it for $N = 2$. The general inductive argument is then obvious and therefore skipped. By e.g. [Var84, Lemma 10.3] it is clear that for any $\nu \in \mathcal{P}(\mathbb{R}^2)$ holds:

$$Ent(\nu|G^2) = Ent(\nu|G^1 \otimes G^1)$$
$$= Ent(p_*^1\nu|G^1) + \int Ent(\nu^{y_1}|G^1)p_*^1\nu(dy_1) \geq \tfrac{\mathcal{T}_2(G^1, p_*^1\nu)}{2} + \int \tfrac{\mathcal{T}_2(G^1, \nu^{y_1})}{2}p_*^1\nu(dy_1),$$

On the other hand, by the bicausal dynamic programming principle (Proposition 5.2) we get

$$\mathcal{T}_{bc,2}(G^2, \nu) = \inf_{\gamma \in \Pi(G^1, p_*^1\nu)} \int \left\{|x_1 - y_1|^2 + \mathcal{T}_2(G^1, \nu^{y_1})\right\} \gamma(dx_1, dy_1)$$
$$= \mathcal{T}_2(G^1, p_*^1\nu) + \int \mathcal{T}_2(G^1, \nu^{y_1})p_*^1\nu(dy_1).$$

$\square$

**Remark 5.6.** *It is clear that equality cannot generally hold in* (CT2)*, since for $N = 1$ we have $\mathcal{T}_2(G, \nu) = \mathcal{T}_{bc,2}(G, \nu)$ and equality in* (T2) *is not usually the case. On the other hand, if $\nu$ is an affine translate of $G^N$, then of course $2Ent(\nu|G^N) = \mathcal{T}_2(G^N, \nu) = \mathcal{T}_{bc,2}(G^N, \nu)$. More generally, there is equality in* (CT2) *provided $\nu^{y_1,...,y_t}(dy_{t+1})$ is an affine translate of $G^1$, for each $t$ and $\nu$-a.e. $y$.*

## 6. Some geometrical aspects of the bicausal case and connections with stochastic programming

At the end of this section we shall prove Theorem 2.9. We start however with a discussion about geometric properties of the space of probability measures endowed with a bicausal Wasserstein distance (equiv. nested distance). We notice first that the Knothe-Rosenblatt rearrangement offers a way to interpolate in a meaningful and non-linear way between stochastic programs. From Remark 5.4 we know that this rearrangement is bicausal and if all conditional probabilities $\mu^{x_1,...,x_n}$ are atomless, then it is induced by a bicausal map which is clearly characterized as the $\mu$-a.s. unique transformation which is left-continuous and increasing w.r.t. lexicographical order and which pushes forward $\mu$ onto $\nu$. Inspired by the concept of displacement interpolation/convexity in optimal transport (as in [Vil03, Chapter 5]) let us define:

**Definition 6.1.** *Let $\pi = \pi(\mu, \nu)$ be the Knothe-Rosenblatt rearrangement as in* (2.8)*. The lexicographical displacement interpolation between $\mu$ and $\nu$ is then defined as the function*

(6.1) $$t \in [0, 1] \mapsto [\mu, \nu]_t := [(x, y) \mapsto (1 - t)x + ty]_*\pi \in \mathcal{P}(\mathbb{R}^N),$$

*If all conditional probabilities $\mu^{x_1,...,x_n}$ are atomless, we also have*

$$[\mu, \nu]_t = ([1 - t]id + tT)_*\mu,$$

where $T = T(\mu, \nu)$ is the Knothe-Rosenblatt map defined in (2.9) and id denotes the identity map in $\mathbb{R}^N$.

Thus we have that $[\mu, \nu]_0 = \mu$ and $[\mu, \nu]_1 = \nu$; therefore the name. We prove now that many stochastic optimization problems are concave along the curves given by (6.1). Recall from (2.13) that the value of a stochastic program with cost $H$ and noise distribution $\eta \in \mathcal{P}(\mathbb{R}^N)$ is given by

$$v(\eta) := \inf_{u_1(), \ldots, u_N()} \int H(x_1, \ldots, x_N, u_1(x_1), u_2(x_1, x_2), \ldots, u_N(x_1, \ldots, x_N))\eta(dx).$$

**Theorem 6.2.** *Let $H : \mathbb{R}^N \times \mathbb{R}^N \to \mathbb{R}$ be bounded from below, concave in the first variable while convex in the second one. Then the functions*

$$t \in [0, 1] \mapsto v([\mu, \nu]_t)$$

*are concave for each $\nu$ and $\mu$ such that all conditional probabilities $\mu^{x_1, \ldots, x_n}$ are atomless, so we may say that $v(\cdot)$ is lexicographic-displacement concave. If further $v(\mu)$ is attained, then*

$$v(\nu) \leq v(\mu) + \inf_{u^* \in \arg\min(v(\mu))} \int \inf_{\xi \in \partial_x H(x, u^*(x))} \xi \cdot [T(x) - x]\mu(dx).$$

The previous result should be seen as a complement to the results of G. Pflug and A. Pichler [Pfl09, PP12], which give conditions under which $|v(\nu) - v(\mu)|$ can be gauged by the value of a bicausal problem between $\mu$ and $\nu$. Indeed, Theorem 6.2 highlights the connection between the most eminent of bicausal maps, i.e. the Knothe-Rosenblatt rearrangement, and multistage stochastic programming. On the other hand, the previous result can be related to [Vil03, Open Problem 5.17], with the caveat that we replaced the role of the Brenier's map by the Knothe-Rosenblatt rearrangement when defining the interpolations.

*Proof of Theorem 6.2.* Take $\mu, \nu, H$ as stated, and $a, b, s, t \in [0, 1]$ with $a + b = 1$. We will show that $v([\mu, \nu]_{at+bs}) \geq av([\mu, \nu]_t) + bv([\mu, \nu]_s)$. We write

$$A_n(x) = (1 - t)x_n + tT^n(x_n; x_1, \ldots, x_{n-1})$$
$$B_n(x) = (1 - s)x_n + sT^n(x_n; x_1, \ldots, x_{n-1})$$
$$C_n(x) = (1 - at - bs)x_n + [at + bs]T^n(x_n; x_1, \ldots, x_{n-1}),$$

with $x = (x_1, \ldots, x_n)$, so we have

$$v([\mu, \nu]_{at+bs}) = \inf_{u_1(), \ldots, u_N()} \int H\Big(C_1(x), \ldots, C_N(x), u_1\big(C_1(x)\big), \ldots, u_N\big(C_1(x), \ldots, C_N(x)\big)\Big)\mu(dx),$$

and by the concavity assumption,

(6.2) $\quad v([\mu, \nu]_{at+bs}) \geq$

$$a\left\{\inf_{u_1(), \ldots, u_N(\cdot)} \int H\Big(A_1(x), \ldots, A_N(x), u_1\big(C_1(x)\big), \ldots, u_N\big(C_1(x), \ldots, C_N(x)\big)\Big)\mu(dx)\right\}$$

$$+ b\left\{\inf_{u_1(), \ldots, u_N(\cdot)} \int H\Big(B_1(x), \ldots, B_N(x), u_1\big(C_1(x)\big), \ldots, u_N\big(C_1(x), \ldots, C_N(x)\big)\Big)\mu(dx)\right\}.$$

If say $t < 1$, we have that $A_1(x) = A_1(x_1)$ is injective (as it is strictly increasing) so defining $\tilde{u}_1(\cdot) = u_1 \circ C_1 \circ A_1^{-1}(\cdot)$ we have $\tilde{u}_1(A_1(x)) = u_1(C_1(x))$. For $n = 2$ we introduce

$$\tilde{u}_2(z_1, z_2) = u_2\big( C_1 \circ A_1^{-1}(z_1) , C_2\big(A_1^{-1}(z_1), A_2^{-1}(z_2|A_1^{-1}(z_1))\big) \big),$$

which is well-defined due to the function $A_2$ being strictly increasing in its second variable, and verify that $\tilde{u}_2(A_1(x), A_2(x)) = u_2(C_1(x), C_2(x))$. Inductively, we obtain easily for each $n \leq N$ a $\tilde{u}_n$ s.t. $\tilde{u}_n(A_1(x), \ldots, A_n(x)) = u_n(C_1(x), \ldots, C_n(x))$. Hence the first term on the r.h.s. of (6.2) is bounded from below by $v([\mu, \nu]_t)$. Using similar arguments for the second term, we see that $v([\mu, \nu]_{at+bs}) \geq av([\mu, \nu]_t) + bv([\mu, \nu]_s)$ holds if $s, t < 1$. The case $s = 1$ or $t = 1$ can be obtained by a limiting and a Komlos-Mazur type argument (with convex combinations; here the convexity assumption is used). A direct argument, inspired by [PP12] and relying in convexity too, is as follows:

$$\int H\big(A_1(x), \ldots, A_N(x), u_1\big(C_1(x)\big), \ldots, u_N\big(C_1(x), \ldots, C_N(x)\big)\big)\mu(dx) \geq$$

$$\mathbb{E}^\mu\left[H\Big(A_1, \ldots, A_N, \mathbb{E}^\mu[u_1(C_1)|A_1, \ldots, A_N], \ldots, \mathbb{E}^\mu[u_N(C_1, \ldots, C_N)|A_1, \ldots, A_N]\Big)\right],$$

by the convexity assumption and Jensen's inequality, so now observing that if e.g. $t = 1$ then $A = T$, which is bicausal from $\mu$ to $\nu$, we get that $\mathbb{E}^\mu[u_n(C_1, \ldots, C_n)|A_1, \ldots, A_N] = \tilde{u}_n(A_1, \ldots, A_n)$ for each $n$ and we conclude as before. The case $s = 1$ is analogous.

For the last statement we obtain by the concavity assumption that

$$H([1-t]x + tT(x), u^*(x)) \leq H(x, u^*(x)) + t\xi(x) \cdot [T(x) - x],$$

where $u^*$ is any optimizer for $v(\mu)$, and $\xi(x)$ is any measurable selection of $x \mapsto \partial_x H(x, u^*(x))$ (the partial superdifferential w.r.t. the first variable). Such a selection exists by [RW98, Theorem 14.56], for which $H$ must be a normal integrand, but this is true by [RW98, Proposition 14.39]. In particular for $t = 1$ and integrating we get

$$\int H(T(x), u^*(x))\mu(dx) \leq v(\mu) + \int \xi(x) \cdot [T(x) - x]\mu(dx).$$

By the same arguments as before (the convexity assumption, the bicausality of $T$ and $T_*\mu = \nu$), the l.h.s. is an upper bound for $v(\nu)$. All in all, if we could find a measurable selector $\xi(\cdot)$ such that

$$\xi(x) \in argmin_{\xi \in \partial_x H(x, u^*(x))} \{\xi \cdot [T(x) - x]\},$$

this would finish the proof. This follows from the measurable maximum theorem [AB06, Theorem 18.19] after observing that $(x, \xi) \mapsto \xi \cdot [T(x) - x]$ is a Carathéodory function and that the correspondence $x \mapsto \partial_x H(x, u^*(x))$ is nonempty compact-valued and weakly measurable (i.e. measurable in the sense of [RW98, Theorem 14.56]). $\square$

Let us now discuss the geometric interpretation that lexicographic displacement interpolation should have. We start from the observation that for costs of the type

$$c_p(x, y) := \sum_{i \leq N} |x_i - y_i|^p,$$

lexicographic displacement interpolation has "constant speed" w.r.t. the associated bicausal / nested distances, namely:

$$(6.3) \qquad \left(\inf_{\gamma \in \Pi_{bc}(\mu, [\mu, \nu]_t)} \int c_p d\gamma\right)^{1/p} = t \left(\inf_{\gamma \in \Pi_{bc}(\mu, \nu)} \int c_p d\gamma\right)^{1/p},$$

whenever $\mu$ has atomless conditional distributions. Indeed, the map $x \mapsto (1-t)x + tT(x)$ by assumption pushes forward $\mu$ into $[\mu, \nu]_t$ and is bicausal, it is further left-continuous and increasing w.r.t. lexicographical order, and so it is the Knothe-Rosenblatt rearrangement of $\mu$ into $[\mu, \nu]_t$ and therefore optimal for the l.h.s. above. Hence

$$\left(\inf_{\gamma \in \Pi_{bc}(\mu, [\mu, \nu]_t)} \int c_p d\gamma\right)^{1/p} = \int \sum_{i \leq N} |x_i - (1-t)x_i - tT^i(x_i; x_1, \ldots, x_{i-1})|^p d\mu$$
$$= t^p \int \sum_{i \leq N} |x_i - T^i(x_i; x_1, \ldots, x_{i-1})|^p d\mu,$$

which is the r.h.s. of (6.3). The argument is of course reminiscent to the Brenier case. This suggests that for the case $p = 2$ one would want to interpret the corresponding bicausal distance as a geodesic length over the space of probability measures (say absolutely continuous ones) when given a differentiable structure and corresponding metric. This is discussed in [Vil03, Chapter 8] for the classical transport case, and no such thing has yet been accomplished for the present bicausal setting. In a way, a first step in this direction was done by Mikami [Mik12], where a Hamilton-Jacobi formulation for the dual of the bicausal problem was derived.

We finally give the missing proof of Theorem 2.9:

*Proof of Theorem 2.9.* By $(EXP)$ and [Vil08, Theorem 22.10 + (22.16)] we have that $\mu^{x_1, \ldots, x_{t-1}}$ satisfies the next $T_1$ functional inequality (for every Borel prob. measure $m$):

$$(6.4) \qquad \mathcal{W}_1(\mu^{x_1, \ldots, x_{t-1}}, m) \leq \frac{\sqrt{2}}{a_t}\left(1 + \log \int e^{a_t x_t^2} \mu^{x_1, \ldots, x_{t-1}}(dx_t)\right)^{1/2} \sqrt{Ent(m|\mu^{x_1, \ldots, x_{t-1}})}$$
$$\leq \frac{\sqrt{2(1+\lambda_t)}}{a_t} \sqrt{Ent(m|\mu^{x_1, \ldots, x_{t-1}})}.$$

Let us denote by $K_{t-1}$ the constant in the r.h.s. above. By triangle inequality and $(LIP)$:

$$\mathcal{W}_1(\mu^{x_1, \ldots, x_{t-1}}, \nu^{y_1, \ldots, y_{t-1}}) \leq \mathcal{W}_1(\mu^{x_1, \ldots, x_{t-1}}, \mu^{y_1, \ldots, y_{t-1}}) + \mathcal{W}_1(\mu^{y_1, \ldots, y_{t-1}}, \nu^{y_1, \ldots, y_{t-1}})$$
$$\leq C \sum_{i < t} |x_i - y_i| + K_{t-1} \sqrt{Ent(\nu^{y_1, \ldots, y_{t-1}}|\mu^{y_1, \ldots, y_{t-1}})},$$

for $\mu \otimes \nu$-a.e. $(x_1, \ldots, x_{t-1}, y_1, \ldots, y_{t-1})$; we are entitled to do this since we may assume w.l.o.g. that $\nu \ll \mu$. Using (Dyn-Pbc) and applying the above computation recursively, one arrives at:

$$\mathcal{W}_{1,bc}(\mu, \nu) \leq \int \nu(dy_1, \ldots, dy_N) \sum_{j < N} K_{N-j-1}(1+C)^j \sqrt{Ent(\nu^{y_1, \ldots, y_{N-j-1}}|\mu^{y_1, \ldots, y_{N-j-1}})},$$

so using Cauchy-Schwartz for the sum and for the integral we get

$$\mathcal{W}_{1,bc}(\mu, \nu) \leq \sqrt{A}\sqrt{B},$$

where

$$B = \int \nu(dy_1, \ldots, dy_N) \sum_{j < N} Ent(\nu^{y_1, \ldots, y_{N-j-1}} | \mu^{y_1, \ldots, y_{N-j-1}}),$$

which by the additive property of the entropy (e.g. [Var84, Lemma 10.3]) equals $Ent(\nu|\mu)$, and where

$$A = 2 \sum_{j < N} (1 + C)^{2j} \frac{(1 + \lambda_{N-j})}{a_{N-j}^2}$$

$\square$

## 7. Counterexamples

Previously we have discussed that the values of the causal and bicausal problems coincide in the case when the starting measure $\mu$ is the product of its marginals and the cost function has a separable structure. The following two examples show that dropping either assumption causes this equality to break down.

**Example 7.1.** *Here $\mu$ is the product of its marginals, but the cost function $c(x_1, x_2, y_1, y_2) = \mathbb{I}_{(x_1,x_2) \neq (y_1,y_2)}$ is non-separable. Take*

$$\mu = 0.16 \, \delta_{(1,1)} + 0.24 \, \delta_{(1,-1)} + 0.24 \, \delta_{(-1,1)} + 0.36 \, \delta_{(-1,-1)},$$

$$\nu = 0.25 \, \delta_{(1,1)} + 0.25 \, \delta_{(1,-1)} + 0.25 \, \delta_{(-1,1)} + 0.25 \, \delta_{(-1,-1)}.$$

*Then an optimal causal transport plan is*

$$\begin{aligned}
\gamma_c = {} & 0.16 \, \gamma((1,1);(1,1)) + 0.24 \, \gamma((1,-1);(1,-1)) + 0.03 \, \gamma((-1,1);(1,1)) \\
& + 0.01 \, \gamma((-1,1);(1,-1)) + 0.2 \ \gamma((-1,1);(-1,1)) + 0.06 \, \gamma((-1,-1);(1,1)) \\
& + 0.05 \ \gamma((-1,-1);(-1,1)) + 0.25 \, \gamma((-1,-1);(-1,-1)),
\end{aligned}$$

*giving the optimal causal value 0.15, and an optimal bicausal one*

$$\begin{aligned}
\gamma_{bc} = {} & 0.16 \, \gamma((1,1);(1,1)) + 0.04 \, \gamma((1,-1);(1,1)) + 0.2 \, \gamma((1,-1);(1,-1)) \\
& + 0.04 \, \gamma((-1,1);(1,1)) + 0.2 \ \gamma((-1,1);(-1,1)) + 0.01 \, \gamma((-1,-1);(1,1)) \\
& + 0.05 \, \gamma((-1,-1);(1,-1)) + 0.05 \, \gamma((-1,-1);(-1,1)) + 0.25 \, \gamma((-1,-1);(-1,-1)).
\end{aligned}$$

*giving the value 0.19.*

**Example 7.2.** *Consider a quadratic-separable cost function and define $\mu$, failing to be the product of its marginals, and $\nu$ as:*

$$\mu = 0.18 \, \delta_{(1,2)} + 0.24 \, \delta_{(1,0)} + 0.18 \, \delta_{(1,-2)} + 0.08 \, \delta_{(-1,2)} + 0.12 \, \delta_{(-1,0)} + 0.2 \, \delta_{(-1,-2)},$$

$$\nu = 0.1 \, \delta_{(1,2)} + 0.26 \, \delta_{(1,-2)} + 0.16 \, \delta_{(-1,2)} + 0.48 \, \delta_{(-1,-2)}.$$

*An optimal causal plan is*

$$\begin{aligned}
\gamma_c = {} & 0.144 \, \gamma((1,0);(1,-2)) + 0.008 \, \gamma((1,0);(-1,2)) + 0.088 \, \gamma((1,0);(-1,-2)) \\
& + 0.1 \, \gamma((1,2);(1,2)) + 0.008 \, \gamma((1,2);(1,-2)) + 0.072 \, \gamma((1,2);(-1,2)) \\
& + 0.108 \, \gamma((1,-2);(1,-2)) + 0.072 \, \gamma((1,-2);(-1,-2)) + 0.12 \ \gamma((-1,0);(-1,-2)) \\
& + 0.08 \, \gamma((-1,2);(-1,2)) + 0.2 \, \gamma((-1,-2);(-1,-2)),
\end{aligned}$$

*giving the value 2.528, and a bicausal optimal one*

$$\begin{aligned}
\gamma_{bc} = {} & 0.144 \, \gamma((1,0);(1,-2)) + 0.096 \, \gamma((1,0);(-1,-2)) + 0.1 \, \gamma((1,2);(1,2)) \\
& + 0.008 \, \gamma((1,2);(1,-2)) + 0.06 \, \gamma((1,2);(-1,2)) + 0.012 \, \gamma((1,2);(-1,-2)) \\
& + 0.108 \, \gamma((1,-2);(1,-2)) + 0.072 \, \gamma((1,-2);(-1,-2)) + 0.02 \, \gamma((-1,0);(-1,2)) \\
& + 0.1 \, \gamma((-1,0);(-1,-2)) + 0.08 \, \gamma((-1,2);(-1,2)) + 0.2 \, \gamma((-1,-2);(-1,-2)),
\end{aligned}$$

*with value 2.72.*

The previous example also shows us that Condition 5.6 and so the monotone regression condition in [Rüs85, Corollary 2], which holds here, is insufficient to guarantee the equality between (Pc) and (Pbc) even for separable costs. Finally, we present an example showing that there may easily be no causal Monge maps even if classical Monge maps do exist.

**Example 7.3.** *In the case, when*

$$\mu = a_1\,\delta_{(1,2)} + a_2\,\delta_{(1,0)} + a_3\,\delta_{(-1,0)} + a_4\,\delta_{(-1,-2)},$$

$$\nu = a_1\,\delta_{(1,2)} + a_3\,\delta_{(1,0)} + a_2\,\delta_{(-1,0)} + a_4\,\delta_{(-1,-2)},$$

*where $a_i, i = 1, \ldots, 4$ are positive numbers that sum up to one and $a_1 \neq a_4$, one can easily observe that there is no causal Monge map pushing forward the former into the latter; i.e. the mass must split. On the other hand, a non-causal Monge map is given by:*

$$T(1,2) = (1,2) \,,\ T(-1,-2) = (-1,-2) \,,\ T(1,0) = (-1,0) \,,\ T(-1,0) = (1,0)$$

$$\gamma^T = a_1\,\gamma((1,2);(1,2)) + a_2\,\gamma((1,0);(-1,0)) + a_3\,\gamma((-1,0);(1,0)) + a_4\,\gamma((-1,-2);(-1,-2)).$$

REFERENCES

[AB06]    C. D. Aliprantis and K. C. Border, *Infinite dimensional analysis*, third ed., Springer, Berlin, 2006, A hitchhiker's guide. MR 2378491 (2008m:46001)

[Bal]     E.J. Balder, *Lectures on young measures*, cahiers de mathematiques de la décision 9517 ed., CEREMADE-Université Paris-Dauphine.

[BG14]    M. Beiglböck and C. Griessler, *An optimality principle with applications in optimal transport and its offsprings*, Submitted, arXiv:1404.7054, 2014.

[BHLP13]  M. Beiglböck, P. Henry-Labordère, and F. Penkner, *Model-independent bounds for option prices—a mass transport approach*, Finance Stoch. **17** (2013), no. 3, 477–501. MR 3066985

[Bog07]   V. I. Bogachev, *Measure theory. Vol. I, II*, Springer-Verlag, Berlin, 2007. MR 2267655 (2008g:28002)

[BS78]    D. P. Bertsekas and S. E. Shreve, *Stochastic optimal control*, Mathematics in Science and Engineering, vol. 139, Academic Press, Inc. [Harcourt Brace Jovanovich, Publishers], New York-London, 1978, The discrete time case. MR 511544 (80d:93081)

[BS11]    M. Beiglböck and W. Schachermayer, *Duality for Borel measurable cost functions*, Trans. Amer. Math. Soc. **363** (2011), no. 8, 4203–4224. MR 2792985 (2012k:49108)

[DGW04]   H. Djellout, A. Guillin, and L. Wu, *Transportation cost-information inequalities and applications to random dynamical systems and diffusions*, Ann. Probab. **32** (2004), no. 3B, 2702–2732. MR 2078555 (2005i:60031)

[DS14]    Y. Dolinsky and H. M. Soner, *Martingale optimal transport and robust hedging in continuous time*, Probab. Theory Related Fields **160** (2014), no. 1-2, 391–427. MR 3256817

[FÜ04]    D. Feyel and A. S. Üstünel, *Monge-Kantorovitch measure transportation and Monge-Ampère equation on Wiener space*, Probab. Theory Related Fields **128** (2004), no. 3, 347–385. MR 2036490 (2004m:60121)

[GHLT14]  A. Galichon, P. Henry-Labordère, and N. Touzi, *A stochastic control approach to no-arbitrage bounds given marginals, with an application to lookback options*, Ann. Appl. Probab. **24** (2014), no. 1, 312–336. MR 3161649

[GL10]    N. Gozlan and C. Léonard, *Transport inequalities. A survey*, Markov Process. Related Fields **16** (2010), no. 4, 635–736. MR 2895086

[GRST15]  N. Gozlan, C. Roberto, P-M. Samson, and P. Tetali, *Kantorovich duality for general transport costs and applications*, Preprint arXiv:1412.7480v4, 2015.

[HN12]    David Hobson and Anthony Neuberger, *Robust bounds for forward start options*, Math. Finance **22** (2012), no. 1, 31–56. MR 2881879

[Kal02]   O. Kallenberg, *Foundations of modern probability*, second ed., Probability and its Applications (New York), Springer-Verlag, New York, 2002. MR 1876169 (2002m:60002)

[Kno57]   H. Knothe, *Contributions to the theory of convex bodies.*, Michigan Math. J. **4** (1957), no. 1, 39–52.

[Kur07]   T.G. Kurtz, *The Yamada-Watanabe-Engelbert theorem for general stochastic equations and inequalities*, Electron. J. Probab **12** (2007), 951–965.

[Las13]   R. Lassalle, *Causal transference plans and their monge-kantorovich problems*, arXiv:1303.6925, 2013.

[Las15]   _____, *Causal transference plans and their monge-kantorovich problems*, Submitted, arXiv:1303.6925.v2, 2015.

[Led01]   M. Ledoux, *The concentration of measure phenomenon*, Mathematical Surveys and Monographs, vol. 89, American Mathematical Society, Providence, RI, 2001. MR 1849347

[McC97]   R. J. McCann, *A convexity principle for interacting gases*, Adv. Math. **128** (1997), no. 1, 153–179. MR 1451422

[Mik12]   T. Mikami, *A characterization of the Knothe-Rosenblatt processes by a convergence result*, SIAM J. Control Optim. **50** (2012), no. 4, 1903–1920. MR 2974723

[Pfl09]   G. Ch. Pflug, *Version-independence and nested distributions in multistage stochastic optimization*, SIAM Journal on Optimization **20** (2009), no. 3, 1406–1420.

[PP12]    G. Ch. Pflug and A. Pichler, *A distance for multistage stochastic optimization models*, SIAM J. Optim. **22** (2012), no. 1, 1–23. MR 2902682

[PP14]    _____, *Multistage stochastic optimization*, Springer Series in Operations Research and Financial Engineering, Springer, Cham, 2014. MR 3288310

[PP15]    _____, *Dynamic generation of scenario trees*, Comput. Optim. Appl. **62** (2015), no. 3, 641–668. MR 3426120

[RS03]    A. P. Ruszczynski and A. Shapiro, *Stochastic programming*, vol. 10, Elsevier Amsterdam, 2003.

[Rüs85]    L. Rüschendorf, *The Wasserstein distance and approximation theorems*, Z. Wahrsch. Verw. Gebiete **70** (1985), no. 1, 117–129. MR 795791 (86m:60004)

[RW98]     R. T. Rockafellar and R. J.-B. Wets, *Variational analysis*, Grundlehren der Mathematischen Wissenschaften [Fundamental Principles of Mathematical Sciences], vol. 317, Springer-Verlag, Berlin, 1998. MR 1491362

[SDR14]    A. Shapiro, D. Dentcheva, and A. Ruszczyński, *Lectures on stochastic programming*, second ed., MOS-SIAM Series on Optimization, Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA; Mathematical Optimization Society, Philadelphia, PA, 2014, Modeling and theory. MR 3242164

[Tal96]    M. Talagrand, *Transportation cost for Gaussian and other product measures*, Geom. Funct. Anal. **6** (1996), no. 3, 587–600. MR 1392331 (97d:60029)

[Var84]    S. R. S. Varadhan, *Large deviations and applications*, CBMS-NSF Regional Conference Series in Applied Mathematics, vol. 46, Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 1984. MR 758258 (86h:60067b)

[Vil03]    C. Villani, *Topics in optimal transportation*, no. 58, American Mathematical Soc., 2003.

[Vil08]    ———, *Optimal transport: old and new*, vol. 338, Springer Science & Business Media, 2008.

[YW71]     T. Yamada and S. Watanabe, *On the uniqueness of solutions of stochastic differential equations*, Journal of Mathematics of Kyoto University **11** (1971), no. 1, 155–167.

[Zae15]    D. Zaev, *On the Monge–Kantorovich Problem with Additional Linear Constraints*, Mat. Zametki **98** (2015), no. 5, 664–683. MR 3438523