

ARTICLE

# The inevitability argument for choice architecture and the evidence-based view

Viktor Ivanković<sup>1,2</sup>  and Andrés Moles<sup>1,2</sup> 

<sup>1</sup>Department of Practical Philosophy, Institute of Philosophy, Zagreb, Croatia and <sup>2</sup>Department of Political Science and Philosophy, Central European University, Vienna, Austria

**Corresponding author:** Andrés Moles; Email: [molesa@ceu.edu](mailto:molesa@ceu.edu)

(Received 20 August 2024; accepted 12 January 2025)

## Abstract

Possibly the most important argument for the permissibility of choice architecture is the inevitability argument (IA), which states that if choice contexts are to inevitably influence individuals in some way, they should be arranged to best promote their welfare. Opponents point out a moral difference between influence from unmodified environments (including environments designed without any thought given to behavioral influence) and environments modified specifically to produce behavioral effects. Only the latter are said to subvert the will of those affected, and thereby raise concerns pertaining to manipulation, mistreatment of rational agency, etc. We argue, however, that if choice architects can reliably predict the behavioral effects of both the unmodified environment and available alternatives, the moral difference between these environments seems insignificant. In such cases, one version of the IA persists. This version establishes the permissibility of choice architecture in circumstances of reliable prediction, but is itself neutral toward available normative directions.

**Keywords:** nudge; choice architecture; inevitability; alien control; reliable prediction

## Introduction

For nearly two decades, ‘nudges’, subtle tweaks in choice environments that predictably steer people’s behavior without restricting their option sets or changing their economic incentives (Thaler and Sunstein, 2008, 6), have been recommended to policy-makers. Inspired by findings in cognitive science about predictable traits in human decision-making – such as a bias toward the present, weakness of the will or an aversion to losses – nudge proponents have been able to come up with a variety of techniques that successfully tap into these heuristics in order to induce welfare-promoting behavior. These include setting up defaults to promote pension savings, designing physical environments to reduce traffic accidents or placing healthy food at eye level in cafeterias to promote healthy eating.

The practice of nudging has been objected to on many grounds – that it is manipulative (Grüne-Yanoff, 2012; Wilkinson, 2013), that it diminishes the control of individuals over their deliberations (Hausman and Welch, 2010), objectionably imposes values upon its targets and disrespects them (White, 2013), fails to treat people as rational agents (Rozeboom, 2020), fails to preserve freedom of choice despite its supposed libertarian credentials (Rebonato, 2014), etc. But according to the founders of nudge, Thaler and Sunstein, many such objections are immaterial, ‘a literal nonstarter’, given that every choice environment will *inevitably* influence choosers in *some* way (2008, 10–11). Putting nudges to use is permissible, the argument goes, because choice architects (and governments employing them) inevitably find themselves in situations in which they must set up arrangements of one kind or another (*ibid.*). The inevitability argument (IA) is, arguably, ‘Thaler and Sunstein’s most important argument for nudging’ (Grill, 2014, 142), from which they draw the following welfarist lesson: if choice contexts need to be arranged in *some* way, then it is best to arrange them so that they make agents better off than they would be in the face of alternative arrangements (Thaler and Sunstein, 2008, 11).<sup>1</sup>

However, some opponents note (e.g., Hausman and Welch, 2010; Grüne-Yanoff, 2012), and more than a few proponents admit (e.g., Blumenthal-Barby, 2013; Grill, 2014; Engelen, 2019), that while influence on choice might indeed be inevitable, there remains a significant moral difference between influence from unmodified environments (including environments designed without any thought given to behavioral influence), and from environments modified specifically to produce behavioral effects. It is only by influences of the latter kind that citizens can be subject to alien control (Schmidt, 2017) and that power can be exerted over them (Hausman and Welch, 2010, 133). In short, while influence may be inevitable, governments can still refrain from becoming nudgers, thus avoiding to enter into morally suspect relationships with their citizens (Grüne-Yanoff, 2012, 639). If this objection stands, then the standard objections against nudging (pertaining to manipulation, lack of respect, imposition of values, etc.) resurface. Proponents have offered little in the way of refuting this objection, despite the supposed importance of IA for nudging.

In this paper, we argue that in at least some cases, one version of the IA persists. While it does not lead straight to Thaler and Sunstein’s welfarist conclusion, it justifies choice architects in interfering with environments that are either unmodified or that were designed without any intention to produce behavioral influence. When choice architects are able to reliably predict the behavioral effects of both the unmodified environment and the available alternative arrangements, they will be permitted, or so we argue, to treat all these options – modified or unmodified, as well as intended or merely foreseen – as on a moral par. In other words, there is insignificant moral difference in such cases between arranging an environment with predictable effects and allowing an unmodified environment with predictable effects, as well as intending a behavioral effect as opposed to merely foreseeing it. The normative position that little or nothing

---

<sup>1</sup> Similar accounts of IA are offered by Thaler and Sunstein (2003), Sunstein (2015), Cohen (2013), Brooks (2013), and Engelen *et al.* (2018).

hinges on whether an environment has been purposely altered or not is supported by what we call the *evidence-based view*.

We proceed as follows. First, we offer some conceptual preliminaries and cast aside versions of IA that have already been rejected in the literature. Second, we offer our own version of IA and ground it in the evidence-based view. Third, we elaborate in more detail why it follows from the availability of evidence and the predictability of contextual effects that choice architects should be morally permitted to modify environments and treat intended as opposed to merely foreseen behavioral effects as on a moral par. Fourth, we face our view with some pressing objections. The final section concludes.

### On what is and isn't inevitable

The debate on IA has been riddled with misnomers and misconceptions. This is primarily due to a lack of conceptual rigor from its main proponents, who have advanced claims about the 'inevitability of nudging' (Thaler and Sunstein, 2008), the 'inevitability of paternalism' (Thaler and Sunstein, 2003; Glaeser, 2006) and 'the inevitability of choice architecture' (Thaler and Sunstein, 2008; Cohen, 2013; Sunstein, 2014, 2015) almost interchangeably. We start by offering some conceptual preliminaries to distinguish between these claims and those that are clearly false. This clarification will also shed light on how we will be using these terms in the paper.

Four terms ought to be clarified – 'paternalism', 'nudging', 'choice architecture/arrangement' and 'choice environment/context'. Let's take 'paternalism' and 'nudging' together. According to Thaler and Sunstein, an intervention qualifies as a nudge only if it is (a) easy and cheap to avoid (preservation of the option set) and (b) makes individuals better off by their own lights (paternalism) (2008, 6, 10). The first condition ensures that the targeted agents need not invest much effort into resisting the influence. The second condition presupposes means paternalism, according to which it is permissible to interfere with the means with which the targeted agents pursue their ends, but not with the ends that they choose for themselves.<sup>2</sup> To count as a nudge, then, a behavioral intervention ought to be easily resistible and means-paternalistic.

Consider now the pairing of 'nudging' and 'choice architecture' (or 'choice arrangement'). We take nudges to be instances of choice architecture that satisfy the aforementioned qualifications (easy resistibility and means paternalism). Choice architecture, more broadly, represents a conscious intervention with predictable behavioral effects, grounded in findings from cognitive science and behavioral economics, but which is not necessarily constrained by conditions like easy resistibility and means paternalism.

---

<sup>2</sup>A behavioral intervention can also be ends-paternalistic, in which case the paternalist interferes with the targeted agents because she believes they have mistaken, confused or irrational ends (Dworkin, 2020). Although Sunstein explicitly endorses a means paternalism, he admits that it will hardly be possible at present to avoid nudges also interfering with ends (2014, 67–68). For a similar point, see Knies (2021). We can cast this issue aside, since it only matters whether it is inevitable for a behavioral intervention to be paternalistic, not what kind of paternalism it fits best.

‘Nudging’ is thus subordinate to ‘choice architecture’ – every instance of nudging is an instance of choice architecture, but not vice versa.

Finally, consider the pairing of ‘choice architecture’ (or ‘choice arrangement’) and ‘choice environment’ (or ‘choice context’). Choice architectures are those choice environments that were consciously modified with an awareness of how this would predictably affect behavior. They thus include an element of *design*. Yet, a choice environment may affect the decision-making of individuals without any conscious intervention whatsoever. The notion simply points to the context dependency of preferences and choices. Once more, the first term is subordinate to the second – every instance of choice architecture is an instance of choice environment, but not vice versa. So, which of these, if any, is inevitable, and does it matter for the permissibility of intervention?

To illustrate this, consider an all-too familiar example from the nudge literature, the arrangement of cafeteria food items. In Thaler and Sunstein’s original example, Carolyn is an expert about all things behavioral. Imagine that, as a director of food services for a large system of schools, she is in a position to instruct the employees of these schools how food items in cafeterias ought to be displayed. Based on the observation that students tend to choose items more visually salient to them, a number of strategies become open to Carolyn – she can influence students to maximize profits, randomize the layout, steer them toward what *she* conceives to be the best option for them, or toward what *they* conceive to be the best option for them, among others (Thaler and Sunstein, 2008, 2).

To intervene paternalistically seems in no way inevitable. This much is obvious from the example itself, in which only two out of the four available strategies are paternalistic. In this vein, Grüne-Yanoff notes that if ‘the government decides that it has no business in improving people’s welfare through its choice-architecture design, then it does not act paternalistically in this regard’ (2012, 639).<sup>3</sup>

If we can avoid paternalism, then it would quickly follow that we can avoid nudging, at least on the narrow definition offered by Thaler and Sunstein that takes all nudging to be paternalistic. But the original definition of nudging might be too narrow. Some of the paradigmatic examples of nudges are not paternalistic in any sense, such as the organ donation default, which clearly benefits someone other than the person being influenced. In fact, it has become commonplace in the literature to attach the ‘nudge’ label to various interventions affecting prosocial and moral behavior (e.g., Guala and Mittone, 2015; Nagatsu, 2015; Capraro *et al.*, 2019).<sup>4</sup> Kelly envisions nudges that promote the principles of Rawlsian justice (2013, 223–225), whereas Moles argues in favor of nudges that facilitate the fulfilment of enforceable duties (2015, 659–660). Still, even granting a broader, non-paternalistic variety of nudges would not make nudging inevitable, for choice architects could still fail to preserve option-sets and thus run afoul of the easy resistibility requirement. For instance, Carolyn could place desserts in a different location altogether, raising the transaction costs for dessert lovers to the

<sup>3</sup> Similar points have been raised by Mitchell (2005), Salvat (2008) and Barton and Grüne-Yanoff (2015).

<sup>4</sup> Some authors omit means paternalism altogether from the definition of nudges (Saghai, 2013; Moles, 2015).

point that the influence is no longer easy to resist (Thaler and Sunstein, 2003, 1184).<sup>5</sup> Nudging is, thus, not inevitable either.

How about the inevitability of influence of choice environments? Regardless of whether environments are arranged with or without an eye to behavioral influence (or are completely unmodified), influence on the behavior of agents will occur. This much seems true, but trivially so. The inevitability of contextual influence is hardly a point of contention between proponents and opponents of employing nudges and choice architecture. Opponents will admit that individuals cannot escape the effects of contextual influence. They may even grant that influence without a designer can be significant to considerations of personal autonomy. But there remains a significant moral difference, they would insist, between the inescapable effects of unmodified environments, and the effects of environments arranged with an astute awareness of how behavior will be affected. And what of Thaler and Sunstein's claim that, since contextual influence is inevitable, environments ought to be arranged in welfare-promoting ways? Opponents would likely insist that this is a question-begging conclusion, to which the fact of contextual inevitability doesn't point. Hence, the inevitability of contextual influence doesn't seem to do much for choice architecture and nudge proponents.

Finally, is choice architecture inevitable? On the conception that we explicate in the next section, it will be inevitable for choice architects to assume control over the behavioral effects on exposed individuals, when they can reliably predict the outcomes of available choice environments. This, we argue, will place unmodified and arranged choice environments on a moral par.

### Inevitability and the evidence-based view

Opponents of choice architecture insist that it is at least *pro tanto* wrong to intervene on choice environments to produce predictable effects. Even if the effects of unintended influences and influences by design equally affect decision-making, influences by design contain an added threat to autonomy, since they might make our actions dependent on the wills of others. We now argue that this suggestion is not as straightforward as it may seem. There are some cases in which a choice architect will be unable to avoid making the actions of others dependent on their wills.

Let's demote Carolyn from the position of higher-up official to a school cafeteria manager. She remains an expert on behavioral influences of all kinds, and successfully predicts not only the effects of arrangements available to her but also the effects when she does not intervene, or the effects of randomized arrangements. But if she can accurately predict the effects of both her action and inaction, and chooses inaction for a particular effect to be produced, is there a moral difference between choosing modified and unmodified environments? With such reliable evidence about influences, are Carolyn's acts and omissions morally distant enough to imply different moral conclusions? We claim that they are not. This gives rise to the evidence-based view, which

---

<sup>5</sup>Objections against the inevitability of nudging are also raised by Rebonato (2014, 371) and Gelfand (2016, 604).

posits the moral proximity of predictable choice environments, regardless of whether they are modified or not.<sup>6</sup>

The case for inevitability is even stronger in environments that are made entirely ‘from scratch’, with no existing default environment. Imagine that a building is being constructed, or that some new regulation requiring a default is being set up. If choice architects like Carolyn can reliably predict the behavioral effects of available arrangements, they inevitably have to pick *some* arrangement knowing what behavioral effects will likely be produced as a result. The case is even stronger here because bringing about behavioral effects is constitutive in bringing the arrangement into existence, without raising concerns about possible differences between action and inaction.

To illustrate how this new kind of inevitability arises, consider the following case:

*Disgruntled customer:* After reading in a magazine about a cafeteria arrangement that uses visual cues to promote healthy eating, a customer recognizes it in Carolyn’s cafeteria. Disgruntled, he faces Carolyn and complains that he doesn’t take kindly to being manipulated into food choices through her behavioral schemes. Carolyn responds, however, that she has a fairly good idea how her customers will be steered *no matter which arrangement* is put into place. She cannot help but pick *some* arrangement, while knowing what kind of behavior will likely be promoted as a result. While she indeed arranges the cafeteria with behavioral effects in mind, she can hardly be at fault for it, since she had to pick *some* arrangement while knowing what behavioral effects will thereby be promoted.

Compare Carolyn to Naïve Bill. Naïve Bill has been in Heuristics School for a week, and only has a fairly good idea about the effects of *one* behavioral arrangement. Imagine now that Bill gets the chance to manage a cafeteria and test his new insight. Unlike the previous case, the disgruntled customer does seem to have a legitimate complaint against Bill, since Bill could’ve avoided putting his insight to the test. This means that Bill’s intervention is susceptible to the charges raised by the opponent of nudges mentioned earlier, pertaining to manipulation, value imposition, etc.

The moral proximity of choice environments available to Carolyn, be they interfered upon or merely allowed, is established because, due to the availability of evidence and predictability of contextual effects, Carolyn cannot in one sense avoid assuming control over the behavioral effects of the choice environment on her patrons. We elaborate this claim further below. For now, we put forward an adjusted version of IA:

*The Inevitability Argument:* Choice architects who can reliably predict the outcomes of available arrangements and the default environment inevitably assume control over the behavioral effects of their designated choice environment on those affected. In such cases, it makes inconsiderable moral difference whether they choose a modified or unmodified choice environment.

---

<sup>6</sup>Alternatively, Carolyn, the higher-up official, could merely have knowledgeable choice architects in her employ (e.g., Carolyns, the school cafeteria managers), whose knowledge of behavioral effects she has at her disposal and to whom she can give orders regarding how choice architectures should be arranged, for inevitability of this kind to obtain.

A few words on what we mean by ‘control’ here. There are two senses in which choice architects can be said to assume control. First, it is one thing to have control over a choice environment itself. In this sense, Carolyn and Bill have control over the cafeteria just the same, in virtue of having the power to modify its aspects, as would any other person in charge regardless of their behavioral expertise. Another, very different sense, is having control over the behavioral effects of that choice environment, the kind of control that we seem to be interested in when we discuss nudging and choice architecture more broadly. The first kind of control – that over a choice environment – is necessary for control over behavioral effects, but not sufficient. This latter kind of control has an epistemic requirement – of knowing, or being able to predict, how modifications to the environment will affect those exposed to it. To have control over the behavioral effects on others, we are required to have a sense of how the intervention upon the environment will steer their behavior. In this sense, Bill only has control over behavioral effects when he puts his one insight from Heuristics School to use. Otherwise, he could change the environment without knowing how it would affect behavior, or that it would affect behavior at all. Neither would amount to assuming control on this second understanding, that over behavioral effects. This understanding, we believe, captures the kind of control that critics of choice architecture usually find disturbing when they raise concerns about, say, manipulation, and it is this kind that we discuss in the remainder of the paper.

Returning to the cafeteria, imagine that the disgruntled customer is unconvinced by Carolyn’s reasonable appeal to the adjusted IA, and gathers like-minded individuals to stage a protest. Carolyn is sacked as a result, and the cafeteria shortly operates with a skeleton staff, which cluelessly and accidentally puts up a nearly identical arrangement to Carolyn’s. Later, James, an equally knowledgeable choice architect as Carolyn is hired. James sees how the cafeteria has been arranged, and happily leaves it untouched, knowing (by hypothesis) what people are likely to choose by virtue of its effects. But James’s passivity seems no different in moral terms to Carolyn’s activity. The disgruntled customer’s grounds for complaint (or lack thereof) seem to be identical against both managers.<sup>7</sup>

Note that the adjusted IA that we defend isn’t biased in favor of welfarism, pro-social behavior, the prevention of harms or any other positive normative account that favors the implementation of some choice environment. It only establishes that the fact that some environment isn’t modified doesn’t count in its favor. Here we distinguish between the content of the influence and its form. Our argument is that, in these relevant contexts, the permissibility of a choice arrangement depends solely on its content.

Imagine that there are two available school cafeterias, *A* and *B*, where the former is not modified by a knowledgeable choice architect, and the latter is. *B* will promote some healthy foods that will otherwise likely be overlooked. Our version of IA simply establishes that, if the choice architect fits our previous description, it will not count in favor of *A* that it is not modified. Of course, the choice architect must later weigh

---

<sup>7</sup>It could be claimed that it is not clear whether James intends or merely foresees the effects of the cafeteria. We turn to this concern in the next section.

reasons regarding which normative direction seems best, given the available options. Perhaps the availability of options will favor a means-paternalistic arrangement (Thaler and Sunstein, 2008), or one imposing minimal costs (Blumenthal-Barby, 2013, 183), or a more ‘natural’ ordering (White, 2010, 217). Additionally, the legitimacy of the choice architect’s decision may be constrained by the kinds of reasons that may be offered in favor of some available choice arrangement; perhaps, such reasons must be public, so that they can be accessible or acceptable to their targets. But these discussions are separate to a prior one on whether employing choice architecture is itself permissible (Grill, 2014, 153). IA in our version would merely deny, for our range of cases, that options are *pro tanto* wrongfully selected if they are arranged to produce predictable effects. It follows from IA that it is permissible to engage in choice architecture. Thus, on our view, it makes little sense to raise complaints about the utilization of choice architecture itself in the circumstances that we describe (on grounds that it is manipulative, or disrespectful), but perfectly legitimate complaints could still be raised about the normative content that many of these arrangements promote.

Some might suspect that the effects of consciously designed choice arrangements on individuals will contain ‘added force’ (i.e., stronger influence), compared to that of their unmodified predecessors. If choice arrangements are more forceful than unmodified environments, then a *pro tanto* reason might reemerge to stick to the latter. However, while it is indeed possible for some arrangements to produce stronger effects than some unmodified contexts, this seems entirely contingent. Nudges are often proposed in order to dull the effects on already triggered heuristics and are supposedly designed to be easily resistible. Imagine that Carolyn’s customers want to follow balanced diets, but in the current setting (one set up by the previous manager with no behavioral aim), they have a hard time resisting their strong temptations for desserts. It would seem that rearranging the cafeteria to promote balanced diets hardly replaces the weaker influence for the stronger. These kinds of cases, where a weaker influence substitutes a stronger, are not the exception in the nudge literature.

We now turn to further explaining what it is about the availability of evidence and the reliability of prediction that brings about this new version of inevitability for the choice architect.

### The importance of being knowledgeable

In the previous section, we argued that the disgruntled cafeteria customer has no justified claim (grounded in manipulation, value imposition, respect, etc.) against knowledgeable choice architects like Carolyn and James, but does have one against a rookie architect like Naïve Bill. The former are absolved, we suggested, because they reliably predict the behavioral effects of available choice environments, making it inevitable for them to assume control over these effects. But what might be the kind of knowledge, or the level of competence, that sets apart Carolyn and James from Bill? Is such a level realistic enough to be relevant for practical considerations? And why does it deflect the charge pressed by the disgruntled customer? In this section, we provide more detail about the kind of evidence that grounds the adjusted IA.

Let’s start from the first question. Responding to the disgruntled customer, Carolyn says that she has a ‘fairly good idea’ about how her customers will be steered no matter

which arrangement is put into place. How might we understand Carolyn's notion of a 'fairly good idea'? At the very least, Carolyn will have a good track record of predicting which available option will be most promoted by some environmental adjustment, and which option will be most impeded. At best, she will be able to predict changes in her customers' general preference ordering that result from the adjustment. To reliably predict this, Carolyn will need to have a good grasp not only of how particular environments trigger heuristics, but what the predominant preferences and values are in a given community. In short, choice architects 'must have a sufficient grasp of the scientific material and a good understanding of how people think in particular situations' (Selinger and Whyte, 2010, 469). However, note that for the purposes of this paper, we are not looking to provide a comprehensive conception of choice-architect expertise. We are only interested in the nature and level of competence needed in some relevant choice context for our version of inevitability to occur. Note that the cafeteria might be the only context in which inevitability obtains for Carolyn. Inevitability within that context may simply be the result of evidence being more easily attainable for Carolyn, or there only being few arrangement options available.<sup>8</sup> This, of course, is not to deny that greater competence among choice architects will typically breed more occurrences of inevitability.<sup>9</sup> As Grill notes, inevitability seems to be 'strengthened by the fact that our knowledge of behavioral psychology is steadily increasing and spreading' (2014, 143). Or as Blumenthal-Barby has recently stated:

once behavioral science helps us gain insight into how choice is affected, intentionality is forced, in a sense. It becomes increasingly difficult for us to maintain that we did not know how various factors in the choice architecture would impact [...] choice. [...] Given that we then have to make a decision about *how* to set things up, we are forced to engage in nudging or shaping choice one way or the other. (2021, 67)

Here's how Carolyn's required level of expertise might be illustrated. Suppose that Carolyn is arranging a cafeteria to promote vegan dishes in a community mostly populated by meat eaters. She judges that there are two prominent spots in which

---

<sup>8</sup>Consider a prison warden who must decide whether the prison walls should be painted a soothing shade of green or remain in a rousing shade of red. Assume that the warden has scientifically valid grounds to believe one option will attenuate the aggression of prison inmates, while the other option will induce it. Assume also that, for whatever reason, the warden doesn't have any other kind of paint available. On our account, it seems inevitable for the warden to have the walls painted in *some* way while knowing what behavioral effects will likely be promoted. Yet it would be erroneous to suggest that this makes the warden a behavioral expert. Instead, inevitability is generated by a short range of available options. The details of this example are inspired by a thought experiment in Douglas (2018) and the research in Schauss (1979).

<sup>9</sup>One might ask whether a sort of 'semi-inevitability' obtains for individuals whose competence lies somewhere in between that of Carolyn and Bill. Imagine that Bridget can utilize various techniques and predict their effects with reasonable accuracy but is oblivious to the effects of other available options. As we see it, inevitability is a threshold concept, so 'semi-inevitability' would not obtain. Simply, it wouldn't be inevitable for Bridget to pick an environment that would place an exposed person under her behavioral control. But greater competence may still have normative significance. The more options there are with predictable effects for Bridget, the more normatively desirable directions there are that she'll be able to pursue via choice architecture. As a result, the requirement for non-control over behavioral effects seems to bear less moral weight, for it would have to trump a greater number of available, desirable aims that could be pursued.

she can display the food, and that these can be ordered by degree of prominence. However, it occurs to Carolyn that placing the vegan dish in the most prominent spot is unlikely to maximize this dish getting picked, since this increases the chance of ‘reactance’ among meat-eating customers, i.e., the likelihood that they will notice that they are being steered and pull in the opposite direction (Sunstein, 2014, 154). Being able to predict this and knowing that she can do more to promote the dish by avoiding consumer reactance, Carolyn places it in the second most prominent spot.

Still, we noted Carolyn only has a ‘good idea’ about how various choice environments will steer her customers. We understand this notion as permitting the occasional mistake, the occasional overlooking of relevant factors on decision-making or the inability to express the prediction in statistical terms. Yet, a rough prediction of this kind still seems sufficient for saying that Carolyn cannot avoid the consideration of behavioral effects in picking a choice arrangement. This is what we will understand as the lowest threshold for Carolyn to make ‘reliable predictions.’

One more important qualifier remains. Carolyn will not be reliably predicting the behavior of each individual patron in any particular instance of choice. She won’t be able to say how the cafeteria arrangement will influence the disgruntled customer specifically, and what food he’ll personally end up picking as a result. To know this, she would have to be intimately familiar not only with his food preferences, but the values that may guide him in making food choices. Instead, Carolyn can make fairly accurate predictions about how changes in the arrangement will shift preferences on the collective level. This is similar to how choice contexts are arranged for gamblers. Architects of gambling contexts can manage the gamblers’ environment so that they don’t quit while they’re ahead, and arrange that over a series of plays, in most cases, gamblers end up with less money than they started with. But that doesn’t mean architects can reliably predict for any particular person that she will pull the lever on a slot machine or stay at the roulette table.

We have noted earlier that the adjusted IA will only apply ‘in some cases’. This is because there is a limited number of cases in which choice architects can be expected to have fairly good ideas about the effects of available options on those exposed. For such cases, evidence will be more easily attainable, arrangement options will be limited, and the effects of influence more stable across cases. It comes as no surprise that defaults and cafeteria arrangements – Thaler and Sunstein’s go-to examples – best fit this description. Inevitability may also occur for physicians in their interactions with patients, as some nudge proponents point out (Brooks, 2013; Cohen, 2013).

Let’s turn now to the second question – why does inevitability deflect the charges of the disgruntled customer? In a nutshell, the knowledgeable choice architect is accountable for the effects of choice environments, regardless of whether she actively brought them about or merely allowed them; since she is able to predict them, she can cause and prevent them. She is accountable for them in much the same way as she would be for the secondary effects of environments that she does not intend but is able to foresee. Since omissions and active interferences are thus placed on a moral par, or so we

suggest, she shouldn't be charged for merely having used choice architecture. Given her knowledge, acting and omitting are sufficiently morally close.<sup>10</sup>

We hint here that there may not be significant normative difference between the expert who reliably foresees an effect and the expert who intends that same effect. Philosophically, this is an altogether different concern from that of the moral significance of acts as opposed to omissions. One might retort that there is an important moral difference between intending a harmful behavioral effect and merely foreseeing it as a side-effect of a choice architecture arranged primarily for, say, aesthetic purposes. The significance of intentions has, indeed, a long pedigree (see Quinn, 1993; Tadros, 2015; for criticisms about the significance of intentions, see Scanlon, 2010). The literature on it is vast, and we cannot do justice to it here. We limit ourselves to surmising that in circumstances of reasonable certainty about side-effects, as is the case of the knowledgeable choice architect who has a fairly good idea how the persons exposed would be affected by different options of choice arrangements, the moral significance of intention as opposed to mere foreseeability is at least downplayed when foreseen effects are ascertained with a fair degree of certainty. Reliable prediction, as we show, similarly downplays the significance of intention as opposed to mere foreseeability as it was established earlier about the significance of acts as opposed to omissions. While these are separate normative points, they largely overlap in the cases we presented here.

The downplaying of the significance of intentions in cases of reliable prediction can be supported by at least two points. First, the observation that foreseeing may at least sometimes be morally on a par with intending is at least partly confirmed by Knob's experiments (2003), which show that merely foreseen harm is often perceived as (being on a par with) intended harm; in other words, being able to predict harm as a secondary effect of one's action is often perceived by test subjects as the same as intending that harm.<sup>11</sup> Second, we should reiterate Blumenthal-Barby's claim that with added insight, 'intentionality is forced, in a sense', and the choice architect is 'forced to engage in [...] shaping choice one way or the other' (2021, 67); this is to say that being able to reliably foresee how an environment affects choice and having to introduce some choice environment with this knowledge at hand is *akin* to intending. There are two ways to understand this claim. The first suggests a moral closeness between intentions and mere foreseeability brought about by new insights in behavioral science, which resembles our own claim. The second understanding points to an epistemic closeness – allowing and merely foreseeing are not always clearly distinguishable (see, e.g., Fitzpatrick, 2006). If Carolyn ignores the foreseeable harms that might befall patrons with hemochromatosis, it's not clear whether she intends them or merely

---

<sup>10</sup>This does not extend to some general view about the moral significance of acts as opposed to omissions. It might be possible to say that a moral difference between acts and omissions persists, but is reduced to the point of being morally inconsiderable in the cases of predictable choice environments we mention here. Or, for those with stronger consequentialist leanings, the difference might fade away entirely. We remain agnostic on this point. Most importantly, the plausibility of IA does not hinge on a particular view about acts and omissions.

<sup>11</sup>Granted, the experiments show that merely foreseen benefits are usually not perceived as (being on a par with) intended benefits.

foresees them. Both understandings support the downplaying of the significance of intended as opposed to merely foreseen effects in the cases at hand.

We will remain agnostic on just *how much* the reliable predictability of foreseen effects downplays the significance of intentions as opposed to mere foreseeability. To us, and for our purposes, it seems at least that in the cases at hand, it doesn't impact concerns of permissibility, and that the burden of proving a significance impacting permissibility would fall on those with intuitions conflicting with ours about the effects of reasonable certainty regarding side-effects and about inevitability.

We finish this section with some claims from the nudge literature that seem to be pointing in a similar argumentative direction to ours. Saghai argues that if secondary behavioral effects may lead to significant costs to some individuals, even if these effects were unintended but merely predictable, the choice architect (or her superior) would be found accountable for bringing the arrangement about (2013, 492). Consider such a case – a cafeteria manager has the option of prominently positioning iron-supplemented food that would benefit the majority of patrons, but could be very harmful for those suffering from hemochromatosis, a rare condition in which the accumulation of iron adversely affects vital organ systems (Salvat, 2008, 11). However, if the cafeteria manager is accountable for allowing unintended, yet predictable adverse secondary effects, the point seems to carry over to considerations of actively bringing about arrangements as opposed to allowing the effects of unmodified environments – the cafeteria manager would be similarly accountable for allowing predictable adverse effects of unmodified choice environments as she would be for the environments she actively brings about.

It might seem odd that while Carolyn is accountable for secondary effects and the effects of unmodified environments, she is somehow absolved of the disgruntled customer's charge. This is because she cannot be at fault for inevitably picking one available choice environment with behavioral effects that are predictable to her. She must answer for the ways in which these environments have been arranged, but it makes little sense that she answers for having arranged them in the first place.<sup>12</sup>

## Objections

In this section, we clarify our position on inevitability further by addressing some objections from the literature.

### *Engaging with choices reflectively*

The first objection is often raised in discussions on nudge inevitability. Specifically, Mitchell has argued, and others have followed suit (Gelfand, 2016, 605; Holm, 2017, 38–39), that contextual influence is inevitable only insofar as 'individuals remain subject to these irrational influences' (Mitchell, 2005, 1251). But individuals can overcome such influences. For instance, 'simply asking people to give reasons for their choices can

---

<sup>12</sup>Although we remain skeptical about the possibility of intentions' significance for permissibility, plausibly, its significance could depend on whether the intentional harm is inflicted as a means to produce a positive effect (for a defense of this interpretation of the doctrine of double effect see Tadros, 2015). None of the examples we discuss have this feature.

reduce the influence of gain/loss framing effects' (ibid., 1256).<sup>13</sup> Indeed, the notion that targeted individuals can rise above heuristic triggers has led some authors to believe that nudging is permissible only insofar as it is in some sense transparent, allowing dissenters to dodge nudges with which they disagree (Schmidt, 2017; Ivanković and Engelen, 2019). In addition, 'boost' proponents have recommended improving people's decision-making competencies so they are able to avoid heuristic triggers (Grüne-Yanoff and Hertwig, 2016). All these proposals might bring into question not only our version of IA, but the inevitability of contextual influence as well.

We do not wish to challenge either the empirical grounding or the moral desirability of these proposals. However, they wouldn't render our considerations of inevitability moot. Imagine that aside from arranging the cafeteria, Carolyn instructs her staff to put up posters at entry points, warning students about the ways in which their food choices can be influenced. She also instructs them to verbally prompt students to think about the reasons for their food choices. Jeff, a high school senior, stops to inspect the posters several days in a row. He takes the time to listen to the cafeteria staff and heeds their advice...for a while. Later, he becomes preoccupied with getting good scores on finals and his social life. The prompts lose their novelty. Eventually, they fall into the background and fail to stir up Jeff's reflective capacities. Once again, the cafeteria arrangement becomes significant for his food choices. Now, the objectors might say that Jeff allows the choice architect to take at least partial control over his food choices. The prompts are giving him a chance to overcome the influence, which he chooses to ignore. But this response would overlook that Jeff may fail to pay attention as a result of his reflection being redirected. Imagine that, instead of becoming preoccupied, he's faced with numerous choice environments where architects are seeking to stir up his reflective capacities. Surely, he cannot engage with all these environments, since his capacity for reflective engagement is a limited resource; in other words, Jeff has a limited mental bandwidth (Mullainathan and Shafir, 2013). While Mitchell might be right about the possibilities of overcoming heuristic triggers, individuals will only be able to do so in a limited number of cases. IA, as we conceive it, will then remain significant beyond an individual's capacity for reflection, i.e., once he uses up all of his cognitive resources. And while individuals might be prompted more often in some types of choice environments than in others, e.g., because some decisions are more and some less weighty, they will expend their resources for reflection into choice environments as they see fit. Any one choice arrangement may be engaged with reflectively by some, and subtly trigger the heuristics of others. Hence, considerations of inevitability may remain significant for any one of these choice arrangements.

### ***Eliminating choice architecture***

Opponents of nudges and choice architecture more broadly might grant to us that the adjusted IA stands, but only if we hold onto our stipulation that a knowledgeable person remains in the role of choice architect. But why stick to such a stipulation? There are at least two conceivable ways, the objection goes, in which we can rid ourselves of the

---

<sup>13</sup>Mitchell draws on two empirical sources to support this view: Miller and Fagley (1991) and Sieck and Yates (1997).

choice architect – by keeping behavioral experts away from becoming choice architects, or by randomizing layouts. While there are separate practical and moral considerations involved for each of the two methods, they successfully prevent inevitability from occurring for choice architects, essentially by eliminating choice architecture. If true, then there isn't a context in which choice architecture is truly inevitable.

In a general sense, the objection holds water. The possibility of taking the choice architect out of the equation somewhat reinforces the original charge that there is a significant moral difference between unmodified environments (including arrangements without any thought given to behavioral influence) and environments modified to produce specific behavioral effects. Triggering heuristics in circumstances of predictable effects is no longer morally unproblematic if we can in fact 'work our way around' inevitability. However, in many cases, the elimination of choice architecture is either inconceivable or deeply undesirable.

Consider first the possibility of keeping behavioral experts away from positions in which inevitability would manifest for them. A proponent of this method would suggest it was right all along to sack Carolyn and wrong to hire James, the two highly knowledgeable experts. To avoid possibly manipulative or value-imposing influence, a principle should be upheld that would deny behavioral experts employment that includes the arrangement of choice environments. But this would be very difficult to accomplish with ever greater expertise. As Grill points out, avoiding choice architecture becomes 'more and more difficult as behavioral insights are disseminated through the population', and although policy-makers can require the consideration of behavioral effects to be ousted from design decisions, 'such requirements will be difficult to monitor' (2014, 143–144).<sup>14</sup>

Now, on a different note, if some individuals with behavioral insights would manage to land these jobs regardless, then it would be more desirable, morally and practically, to have people as competent as Carolyn to be choice architects than Naïve Bills. In other words, there are good reasons to suggest that if having some behavioral expertise on board is hardly avoidable, then it's better to have as much of it as possible. On the one hand, as we've established in this paper, Carolyn is more likely to find herself in circumstances of inevitability, absolving her from the charges raised by the disgruntled customer. On the other hand, some argue that individuals who find themselves in the role of choice architect, and who are aware of ubiquitous influence on behavior, must arrange choices responsibly (Hansen and Jespersen, 2013, 23), while others suggest this can be done only at a level of competence that can be trusted (Selinger and Whyte, 2010, 462). Blumenthal-Barby suggests that arranging the choice environment 'should be based on data about satisfaction and happiness levels across various outcomes' (2013, 196).

But even if it were insisted that individuals with behavioral insights could conceivably be kept away from choice architecture positions, thereby overcoming IA, it would still be far from conclusive that all the possible benefits of choice architecture should be entirely forgone to overcome the threats of manipulation or value imposition. Nudges

<sup>14</sup> Additionally, recall that the epistemic bar we set for inevitability to occur is often quite low. In some circumstances, evidence will be more attainable, and/or available arrangement options will be few. Even laypersons may find themselves grasping behavioral effects in such circumstances after a while.

are often proposed in situations in which ‘people’s psychological set-up predictably leads them astray—failing to live up to [...] their own professed values and ideals’ (Engelen *et al.*, 2018, 351). Giving up on choice architecture will often entail leaving individuals at the mercy of their own psychological deficiencies, or, in some cases, to the exploitation of their heuristics by profit maximizers.<sup>15</sup> And even if their heuristics are not strictly used to harm them, individuals would still be missing out on guidance in important areas, including health, wealth and safety.

Alternatively, some authors believe we can avoid inevitability by randomizing the layout. Randomizing is said to offer choice contexts in which no person is ‘under anyone else’s control’ (Wilkinson, 2013, 343). The merit of a randomized layout is neutrality in the sense that no option is being consciously steered by a designer.

While randomization seems conceivable in some circumstances, it might be more difficult to imagine in others. If a landscape architect was designing a park, aware of the behavioral effects that the available design options might bring about, it seems hard to envisage what ‘randomizing a park layout’ would entail. Or consider a doctor aware of the various framing effects at work when presenting a diagnosis or therapy options – once the doctor’s help has been sought, it hardly seems imaginable that he can randomize the frame (Cohen, 2013, 9). Finally, to attend to our favored example here, it’s not at all clear how we should envision a randomized cafeteria.

In many other circumstances, randomization may just seem ‘silly’ (Grill, 2014, 143) – perhaps overcoming inevitability through randomization would be conceivable (and thus could eliminate choice architecture), but hardly favorable. If an environment were truly randomized, then the choice architect would lack veto control over some randomized layout, and this, at worst, would threaten the authorization of some environments with predictably harmful effects (like the iron-supplemented food arrangement is harmful to people with hemochromatosis). Such harmful arrangements could be disqualified, but this would require at least some degree of interference from the choice architect, thereby ruling out randomization in the strongest sense. A weaker kind of randomization might entail picking one of the remaining arrangements at random or randomly alternating between arrangements. However, if we’ve already ‘dirtyed our hands’ with choice architecture to eliminate harmful arrangements, why not also rule out arrangements completely lacking in benefit? But it might be questioned at this point, or even at the previous, whether such randomization truly eliminates choice architecture.

Hence, it’s not at all obvious that the gains from avoiding inevitability in these described ways would outweigh the very significant costs. The objection does, however, force us to concede that eliminating the choice architect is in some cases conceivable, regardless of how normatively preposterous eliminating choice architecture may seem in these cases. Be that as it may, our adjusted IA remains relevant at least for a number of cases where elimination is inconceivable.<sup>16</sup>

<sup>15</sup>For a comprehensive account of the objectionability of behavioral influences in the market, see Ivanković and Engelen (2024).

<sup>16</sup>Future developments of artificial intelligence might eliminate the role of choice architects, but this prospect still seems to be in the distant future.

## Conclusion

We have argued in this paper that, in conditions in which choice architects can make reliable predictions about the behavioral effects of available choice environments, it becomes inevitable for the architect to pick *some* choice environment that generates *some* predictable behavioral effect. In such cases, our argument goes, there is inconsiderable moral difference between picking an unmodified and modified choice environment, as well as between intending and merely foreseeing behavioral effects. Our version of the IA, grounded on the evidence-based view, is itself neutral regarding normative content – it says nothing about the direction in which choice environments should steer. The argument only shows that, in many cases, because of inevitability, the debate about form should be altogether skipped, and turn to content.

**Acknowledgements.** The paper is in part inspired by sections from Ivanković’s doctoral dissertation *The Liberal Politics of Behavioral Enhancement* (Central European University, 2019). We would like to thank Tom Douglas, Tom Parr, Kalle Grill, Anamarija Komesarović, Lovro Savić, Aleksandar Simić, Maximilian Kiener, and two anonymous reviewers for their insightful comments and suggestions, as well as the audiences at the ‘3<sup>rd</sup> Polemo Conference 2021’ (Central European University, June 2021), ‘Society of Applied Philosophy Annual Conference’ (SAP, July 2021), and ‘Influethics seminar’ (Université Catholique de Lille, September 2021).

**Funding statement.** Ivanković was supported by the project Ethics and Social Challenges (EDI) at the Institute of Philosophy, reviewed by the Ministry of Science and Education of the Republic of Croatia and financed through the National Recovery and Resilience Plan 2021–2026 of the European Union – NextGenerationEU.

**Competing interests.** The authors declare that they have no competing interests.

## References

- Barton, A. and T. Grüne-Yanoff (2015), ‘From libertarian paternalism to nudging—and beyond’, *Review of Philosophy and Psychology*, **6**(3): 341–359.
- Blumenthal-Barby, J. S. (2013), Choice Architecture: a mechanism for improving decisions while preserving liberty?, in C. Coons and M. Weber (eds), *Paternalism: Theory and Practice*, New York: Cambridge University Press, 178–196.
- Blumenthal-Barby, J. S. (2021), *Good Ethics and Bad Choices: The Relevance of Behavioral Economics for Medical Ethics*, Cambridge, MA: The MIT Press.
- Brooks, T. (2013), ‘Should we nudge informed consent’, *The American Journal of Bioethics*, **13**(6): 22–23.
- Capraro, V., G. Jagfeld, R. Klein, M. Mul, and I. van de Pol (2019), ‘Increasing altruistic and cooperative behaviour with simple moral nudges’, *Scientific Reports*, **9**(1): 11880.
- Cohen, S. (2013), ‘Nudging and informed consent’, *American Journal of Bioethics*, **13**(6): 3–11.
- Douglas, T. (2018), Neural and environmental modulation of motivation, in D. Birks and T. Douglas (eds), *Treatment for Crime: Philosophical Essays on Neurointerventions in Criminal Justice*, Oxford: Oxford University Press, 208–224.
- Dworkin, G. (2020), Paternalism, *The Stanford Encyclopedia of Philosophy*, in E. N. Zalta (ed.), <https://plato.stanford.edu/entries/paternalism/> (accessed on February 11, 2022).
- Engelen, B. (2019), ‘Ethical criteria for health-promoting nudges: a case-by-case analysis’, *American Journal of Bioethics*, **19**(5): 48–59.
- Engelen, B., A. Thomas, A. Archer and N. van de Ven (2018), ‘Exemplars and nudges: combining two strategies for moral education’, *Journal of Moral Education*, **47**(3): 346–365.
- Fitzpatrick, W. J. (2006), ‘The intend/foresee distinction and the problem of “closeness”’, *Philosophical Studies*, **128**(3): 585–617.

- Gelfand, S. D. (2016), 'The meta-nudge – a response to the claim that the use of nudges during the informed consent process is unavoidable', *Bioethics*, **30**(8): 601–608.
- Glaeser, E. L. (2006), 'Paternalism and psychology', *The University of Chicago Law Review*, **73**(1): 133–156.
- Grill, K. (2014), 'Expanding the nudge: designing choice contexts and choice contents', *Rationality, Markets and Morals*, **5**(90): 139–162.
- Grüne-Yanoff, T. (2012), 'Old wine in new casks: libertarian paternalism still violates liberal principles', *Social Choice and Welfare*, **38**(4): 635–645.
- Grüne-Yanoff, T. and R. Hertwig (2016), 'Nudge versus boost: how coherent are policy and theory?', *Minds and Machines*, **26**(1-2): 149–183.
- Guala, F. and L. Mittone (2015), 'A political justification of nudging', *Review of Philosophy and Psychology*, **6**(3): 385–395.
- Hansen, P. G. and A. M. Jespersen (2013), 'Nudge and the manipulation of choice: a framework for the responsible use of the nudge approach to behaviour change in public policy', *European Journal of Risk Regulation*, **4**(1): 3–28.
- Hausman, D. M. and B. Welch (2010), 'Debate: to nudge or not to nudge', *Journal of Political Philosophy*, **18**(1):123–136.
- Holm, S. (2017), 'Authenticity, best interest, and clinical nudging', *Hastings Center Report*, **47**(2): 38–40.
- Ivanković, V. and B. Engelen (2019), 'Nudging, transparency, and watchfulness', *Social Theory and Practice*, **45**(1): 43–73.
- Ivanković, V. and B. Engelen (2024), 'Market nudges and autonomy', *Economics and Philosophy*, **40**(1): 138–165.
- Kelly, J. T. (2013), Libertarian paternalism, utilitarianism, and justice, in C. Coons, and M. Weber (eds), *Paternalism: Theory and Practice*, New York: Cambridge University Press, 216–230.
- Kniess, J. (2021), 'Libertarian paternalism and the problem of preference architecture', *British Journal of Political Science*, **52**(2): 921–933.
- Knobe, J. (2003), 'Intentional action and side-effects in ordinary language', *Analysis*, **63**(3): 190–194.
- Miller, P. M. and N. S. Fagley (1991), 'The effects of framing, problem variations, and providing rationale on choice', *Personality and Social Psychology Bulletin*, **17**(5): 517–522.
- Mitchell, G. (2005), 'Libertarian paternalism is an oxymoron', *Northwestern University Law Review*, **99**(3): 1245–1278.
- Moles, A. (2015), 'Nudging for liberals', *Social Theory and Practice*, **41**(4): 644–667.
- Mullainathan, S. and E. Shafir (2013), *Scarcity: Why Having Too Little Means So Much*, New York: Times Books, Henry Holt and Company.
- Nagatsu, M. (2015), 'Social nudges: their mechanisms and justification', *Review of Philosophy and Psychology*, **6**(3): 481–494.
- Quinn, W. (1993), *Morality and Action*, Cambridge: Cambridge University Press.
- Rebonato, R. (2014), 'A critical assessment of libertarian paternalism', *Journal of Consumer Policy*, **37**(3): 357–396.
- Rozeboom, G. J. (2020), 'Nudging for rationality and self-governance', *Ethics*, **131**(1): 107–121.
- Saghai, Y. (2013), 'Salvaging the concept of nudge', *Journal of Medical Ethics*, **39**(8): 487–493.
- Salvat, C. (2008), Is libertarian paternalism an oxymoron?: a comment on Sunstein and Thaler. hal-00336528.
- Scanlon, T. M. (2010), *Moral Dimensions*, Cambridge: Harvard University Press.
- Schauss, A. G. (1979), 'Tranquilizing effect of color reduces aggressive behavior and potential violence', *Journal of Orthomolecular Psychiatry*, **8**(4): 218–221.
- Schmidt, A. T. (2017), 'The power to nudge', *American Political Science Review*, **111**(2): 404–417.
- Selinger, E. and K. P. Whyte (2010), 'Competence and Trust in Choice Architecture', *Technology & Policy*, **23**(3): 461–482.
- Sieck, W. and F. J. Yates (1997), 'Exposition effects on decision making: choice and confidence in choice', *Organizational Behavior and Human Decision Processes*, **70**(3): 207–219.
- Sunstein, C. S. (2014), *Why Nudge? The Politics of Libertarian Paternalism*, New Haven & London: Yale University Press.
- Sunstein, C. S. (2015), 'Nudges, agency, and abstraction: a reply to critics', *Review of Philosophy and Psychology*, **6**(3): 511–529.
- Tadros, V. (2015), 'Wrongful intentions without closeness', *Philosophy and Public Affairs*, **43**(1): 52–74.

- Thaler, R. H. and C. S. Sunstein (2003), Libertarian paternalism is not an oxymoron, *The University of Chicago Law Review*, **70**(4): 1159–1202.
- Thaler, R. H. and C. S. Sunstein (2008), *Nudge: Improving Decisions about Health, Wealth, and Happiness*, New Haven: Yale University Press.
- White, M. D. (2010), 'Behavioral law and economics: the assault on consent, will, and dignity', in C. Fodor, G. F. Gaus and J. Lamont (eds), *Essays on Philosophy, Politics & Economics*, Stanford: Stanford University Press, 201–224.
- White, M. D. (2013), *The Manipulation of Choice: Ethics and Libertarian Paternalism*, New York: Palgrave Macmillan.
- Wilkinson, T. M. (2013), Nudging and manipulation, *Political Studies*, **61**(2): 341–355.