

CONTINUOUS NEWTON-LIKE METHODS FEATURING INERTIA AND VARIABLE MASS*

CAMILLE CASTERA[†], HEDY ATTOUCH[‡], JALAL FADILI[§], AND PETER OCHS[¶]

Abstract. We introduce a new dynamical system, at the interface between second-order dynamics with inertia and Newton’s method. This system extends the class of inertial Newton-like dynamics by featuring a time-dependent parameter in front of the acceleration, called *variable mass*. For strongly convex optimization, we provide guarantees on how the Newtonian and inertial behaviors of the system can be non-asymptotically controlled by means of this variable mass. A connection with the Levenberg–Marquardt (or regularized Newton’s) method is also made. We then show the effect of the variable mass on the asymptotic rate of convergence of the dynamics, and in particular, how it can turn the latter into an accelerated Newton method. We provide numerical experiments supporting our findings. This work represents a significant step towards designing new algorithms that benefit from the best of both first- and second-order optimization methods.

Key words. Optimization, Dynamical Systems, Newton’s Methods, Convex Optimization, Differential Equations

MSC codes. 37N40, 46N10, 49M15, 65B99, 65K10

Dedicated to the memory of Hedy Attouch, outstanding mathematician and beloved collaborator.

1. Introduction.

1.1. Problem Statement. A major challenge in modern unconstrained convex optimization consists in building fast algorithms while maintaining low computational cost and memory footprint. This plays a central role in many key applications such as large-scale machine learning problems or data processing. The problems we are aiming to study are of the form

$$\min_{x \in \mathbb{R}^n} f(x).$$

Large values of n demand for algorithms at the interface of first- and second-order optimization. Limited computational capabilities explain why gradient-based (first-order) algorithms remain prominent in practice. Unfortunately, they often require many iterations, which is true even for the provably best algorithms for certain classes of optimization problems; for example that of convex and strongly convex functions with Lipschitz continuous gradient [43, 38, 39]. On the other hand, algorithms using second-order information (the Hessian of f)—with Newton’s method as prototype—adapt locally to the geometry of the objective, allowing them to progress much faster towards a solution. However, each iteration comes with high computational and memory costs, which highlights a challenging trade-off. It is therefore essential to develop algorithms that take the best of both worlds. Several quasi-Newton algorithms partly address this issue, for example BFGS methods [22, 28, 30, 45, 36], yet, in very large-scale applications, first-order algorithms often remain the preferred choice.

* Submitted to the editors 14/10/2023.

Funding: C. Castera, J. Fadili and P. Ochs are supported by the ANR-DFG joint project TRINOM-DS under number ANR-20-CE92-0037-01.

[†]University of Tübingen, Germany (camille.castera@protonmail.com).

[‡]IMAG, Université Montpellier, CNRS, France.

[§]ENSICAEN, Normandie Université, CNRS, GREYC, France.

[¶]Saarland University, Germany.

In order to reach a new level of efficiency, deep insights into the mechanism and relations between algorithms are required. To that aim, an insightful approach is to see optimization algorithms as discretization of ordinary differential equations (ODEs): for small-enough step-sizes, iterates can be modeled by a continuous-time trajectory [37, 16]. Obtaining a fast algorithm following this strategy depends on two ingredients: choosing an ODE for which rapid convergence to a solution can be proved, and discretizing it with an appropriate scheme that preserves the favorable properties of the ODE.

Both steps are highly challenging, our work focuses on the ODE matter. We study the following second-order dynamical system in a general setting:

$$(VM-DIN-AVD) \quad \varepsilon(t)\ddot{x}(t) + \alpha(t)\dot{x}(t) + \beta\nabla^2 f(x(t))\dot{x}(t) + \nabla f(x(t)) = 0, \quad t \geq 0,$$

where $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is a smooth function, with gradient ∇f and Hessian $\nabla^2 f$ defined on \mathbb{R}^n equipped with scalar product $\langle \cdot, \cdot \rangle$, and induced norm $\|\cdot\|$. The system is controlled by two functions $\varepsilon, \alpha: \mathbb{R}_+ \rightarrow \mathbb{R}_+$ (where $\mathbb{R}_+ = [0, +\infty[$) and a parameter $\beta > 0$ that define the type of dynamics that drives the trajectory (or solution) $x: \mathbb{R}_+ \rightarrow \mathbb{R}^n$, whose first- and second-order derivatives are denoted \dot{x} and \ddot{x} respectively. We call the above dynamics (VM-DIN-AVD), which stands for “Variable Mass Dynamical Inertial Newton-like system with Asymptotically Vanishing Damping” since it generalizes a broad class of ODEs whose original member is DIN [2], where ε and α were constant. DIN was then extended to the case of non-constant *asymptotically vanishing dampings* (AVD) α [12]. In this work we introduce the non-constant parameter ε called *variable mass* (VM) in front of the acceleration \ddot{x} , in the same way that α is called (viscous) *damping* by analogy with classical mechanics. A key feature of these ODEs, that positions them at the interface of first- and second-order optimization, is that they possess equivalent forms involving only ∇f but not $\nabla^2 f$, significantly reducing computational costs, hence enabling the design of practical algorithms, see e.g., [24, 8, 25]. The key idea behind this is the relation $\nabla^2 f(x(t))\dot{x}(t) = \frac{d}{dt}\nabla f(x(t))$, see Section 2 for an equivalent formulation of (VM-DIN-AVD) exploiting this.

This paper emphasizes the relation between (VM-DIN-AVD) and special cases. Indeed, taking $\varepsilon = \alpha = 0$, one obtains¹ the Continuous Newton (CN) method [29]

$$(CN) \quad \beta\nabla^2 f(x_N(t))\dot{x}_N(t) + \nabla f(x_N(t)) = 0, \quad t \geq 0,$$

known notably for being invariant to affine transformations and yielding fast vanishing of the gradient (see Section 3). In fact, this observation shows that (VM-DIN-AVD) is a singular perturbation of (CN), which also justifies the terminology “Newton-like” in DIN. When $\alpha \neq 0$ but $\varepsilon = 0$, we recover the Levenberg–Marquardt (LM) method,

$$(LM) \quad \alpha(t)\dot{x}_{LM}(t) + \beta\nabla^2 f(x_{LM}(t))\dot{x}_{LM}(t) + \nabla f(x_{LM}(t)) = 0, \quad t \geq 0,$$

also known as regularized Newton method since it stabilizes (CN). In the rest of the paper, the solutions of (CN) and (LM) will always be denoted by x_N and x_{LM} respectively. Since the introduction of DIN, it is known (see [2]) that for $\alpha = 0$, $\beta = 1$, and ε constant and small, (VM-DIN-AVD) is a “perturbed” Newton method since the distance between the solutions of (VM-DIN-AVD) and (CN) is at most proportional to $\sqrt{\varepsilon}$ at all times. Yet, despite the benefits of this class of ODEs, such as stabilization properties [12, 8], no improvement² of the rate of convergence (in values) has been shown compared to inertial first-order dynamics [43, 47]. This raises the question:

¹CN is usually considered with $\beta = 1$, we put β in the system to ease the discussions below.

²DIN-like systems were thought to yield faster vanishing of the gradient compared to inertial first-order dynamics, until recently [9].

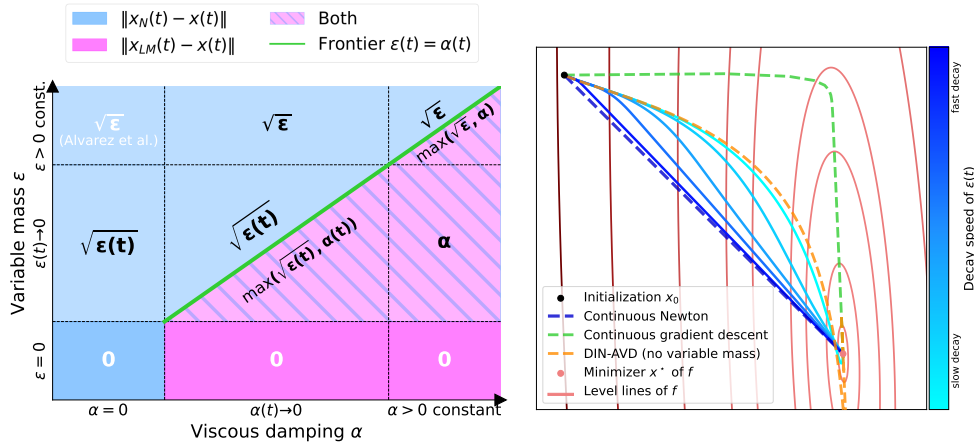


FIG. 1. *Left: phase diagram on distances from (VM-DIN-AVD) to (CN) and (LM) (see Section 3). The color of each patch indicates which distance is considered, and the scaling of a corresponding upper-bound on this distance is written (in white for prior work, in black for our contributions). The green line separates the cases $\varepsilon \geq \alpha$ (above) and $\varepsilon \leq \alpha$ (below). Right: 2D illustration of the trajectories of (VM-DIN-AVD) for several choices of ε on a quadratic function. Fast-vanishing $\varepsilon(t)$ (dark-blue solid curves) bring solutions of (VM-DIN-AVD) close to that of (CN), making them, more robust to bad conditioning compared to first-order dynamics (e.g., gradient descent).*

“are these ODEs really of Newton type?”,

which is crucial in view of designing faster algorithms from them.

1.2. Main Contributions. We show that the answer to this question is partially positive, and closely related to the choices of ε and α . We provide general results on the role played by these two control parameters and how they can be chosen to control (VM-DIN-AVD), and make it close to (CN) *for all time*, as illustrated on the right-hand side of Figure 1, but also to obtain fast convergence. This represents a first step towards building new fast practical algorithms. Our main contributions are the following:

- We provide a first-order equivalent formulation of (VM-DIN-AVD), and show the existence and uniqueness of the solutions of (VM-DIN-AVD) under mild assumptions.
- We generalize the perturbed Newtonian property discussed above to non-constant and possibly vanishing variable masses ε , and “not too large” positive dampings α , and derive bounds that (formally) take the form $\|x(t) - x_N(t)\| = O(\sqrt{\varepsilon(t)})$. We then extend these results to larger α and make the connection between (VM-DIN-AVD) and (LM). This contribution is summarized in the phase diagram of Figure 1.
- Using quadratic functions as a model for strongly convex functions, we shed light on techniques to efficiently approximate solutions of (VM-DIN-AVD). We then show how ε and α affect the speed of convergence. Depending on their setting, the solutions of (VM-DIN-AVD) may either converge as fast as that of (CN), *faster*, or rather have a (LM) nature, as summarized in Table 1.
- We provide numerical experiments supporting our theoretical findings.

1.3. Related work. The importance and challenges of finding systems that approximate and preserve the benefits of (CN) were highlighted by [3]. The system (VM-DIN-AVD) belongs to the class of inertial systems with viscous and geometric (“Hessian-driven”) dampings, initially introduced with constant $\varepsilon = 1$ and constant

TABLE 1
Informal summary of Section 4. Comparison of (VM-DIN-AVD) with other dynamics

Parameters of (VM-DIN-AVD)		Speed of convergence	
Dominant parameter	Integrability in $+\infty$	w.r.t. (CN)	w.r.t. (LM)
variable mass ε	yes	as fast	as fast
	no	faster	faster
viscous damping α	yes	as fast	only depends on ε
	no	slower	

α in [2] and called DIN (for Dynamical Inertial Newton-like system). Except in a few cases [24, 23], most of the follow-up work then considered extensions of DIN with non-constant AVD α , with in particular the DIN-AVD system with $\alpha(t) = \alpha_0/t$ as introduced in [12]. The reason for this popular choice for α is its link with Nesterov’s method [47]. Non-constant choices for β have been considered [8, 1, 34, 5]. We keep it constant, and rather focus on non-constant ε , unlike prior work that used constant $\varepsilon = 1$. The mass ε was only considered in the original work [2], but only for fixed ε , $\beta = 1$ and constant $\alpha = 0$. VM-DIN-AVD is however related to the IGS system considered in [5] as it is actually equivalent to the latter after dividing both sides of (VM-DIN-AVD) by $\varepsilon(t)$. Our approach—which consists in studying the connections with other second-order dynamics as ε vanishes asymptotically—is however different from the one followed in [5], which is of independent interest. The literature on DIN is rich, let us mention further connections with Nesterov’s method [46, 1], extensions with Tikhonov regularization [19] and closed-loop damping [6, 34]. The non-smooth and possibly non-convex cases have been considered in [10, 11, 24]. Finally, avoidance of strict saddle points in smooth non-convex optimization has been shown in [23].

The influence of the damping α on the (LM) dynamics has been studied in [14, 13]. Interestingly, the conditions enforced on α in these papers (formally a sub-exponential decay) are very similar to those we make on ε and α for (VM-DIN-AVD) (see Assumptions 2 and 3). The (LM) dynamics is also related to the system in [4].

Regarding the second part of our analysis, which deals with the case where f is quadratic, a recent work [8] provided closed-form solutions to (VM-DIN-AVD) for $\varepsilon \equiv 1$ and special choices of α . Our work rather deals with approximate solutions which allows considering a wide class of functions ε and α . We rely on the Liouville–Green (LG) method [35, 31] presented in Section 4. Generalizations of LG are also often referred to as WKB methods [50, 33, 21] and seem to be mostly used in physics so far. To the best of the authors’ knowledge, the current work seem to be one of the first to use the LG method in optimization, and the first for DIN-like systems.

1.4. Organization. The paper is organized as follows. We discuss the existence of solutions in Section 2. Our main results, from a non-asymptotic control perspective, are then presented in Section 3. An analysis of the role played asymptotically by ε and α is then carried out on quadratic functions in Section 4. Finally, numerical experiments are presented in Section 5, and some conclusions are then drawn.

2. Existence and Uniqueness of Solutions. We always assume the following on the system (VM-DIN-AVD) and the functions f, ε, α defined in the introduction.

ASSUMPTION 1. • $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex, twice continuously differentiable, coercive, and strongly convex on bounded subsets of \mathbb{R}^n . We fix $x_0, \dot{x}_0 \in$

\mathbb{R}^n , such that, unless stated otherwise, (VM-DIN-AVD) has initial condition $(x(0), \dot{x}(0)) = (x_0, \dot{x}_0)$, and (CN) and (LM) have initial conditions $x_N(0) = x_{LM}(0) = x_0$.

- $\alpha : \mathbb{R}_+ \rightarrow \mathbb{R}$ is non-increasing, non-negative and differentiable. $\varepsilon : \mathbb{R}_+ \rightarrow \mathbb{R}$ is non-increasing, positive, and twice differentiable with bounded second derivative. We fix initial values: $\varepsilon(0) = \varepsilon_0 > 0$, $\varepsilon'(0) = \varepsilon'_0 \leq 0$ and $\alpha(0) = \alpha_0 \geq 0$.

THEOREM 2.1. *Under Assumption 1 there exists a unique global solution $x : \mathbb{R}_+ \rightarrow \mathbb{R}^n$ to (VM-DIN-AVD).*

The proof relies on the Cauchy–Lipschitz Theorem, we sketch the main elements. We reformulate (VM-DIN-AVD) into a first-order (in time) system by introducing an auxiliary variable $y : \mathbb{R}_+ \rightarrow \mathbb{R}^n$. Notably, our reformulation does not involve $\nabla^2 f$, in the same fashion as [2, 12]. For all t , defining $\nu(t) = \alpha(t) - \varepsilon'(t) - \frac{1}{\beta}\varepsilon(t)$, we show in Appendix A that (VM-DIN-AVD) is equivalent to

$$(gVM-DIN-AVD) \quad \begin{cases} \varepsilon(t)\dot{x}(t) + \beta\nabla f(x(t)) + \nu(t)x(t) + y(t) & = 0 \\ \dot{y}(t) + \nu'(t)x(t) + \frac{\nu(t)}{\beta}x(t) + \frac{1}{\beta}y(t) & = 0 \end{cases}$$

with initial conditions $(x(0), y(0)) = \left(x_0, -\varepsilon_0\dot{x}_0 - \beta\nabla f(x_0) - (\alpha_0 - \varepsilon'_0 - \frac{1}{\beta}\varepsilon_0)x_0\right)$.

One can notice that in the special case where ε is taken constant and equal to 1 (that is when (VM-DIN-AVD) is simply the DIN-AVD system [12]), we recover the same first-order formulation as that in [12]. For all $t \geq 0$ and $(u, v) \in \mathbb{R}^n \times \mathbb{R}^n$, define

the mapping $G(t, (u, v)) = \begin{pmatrix} \frac{1}{\varepsilon(t)}(-\beta\nabla f(u) - \nu(t)u - v) \\ -\nu'(t)u - \frac{\nu(t)}{\beta}u - \frac{1}{\beta}v \end{pmatrix}$, so that (gVM-DIN-AVD)

rewrites $(\dot{x}(t), \dot{y}(t)) = G(t, (x(t), y(t)))$ for all $t \geq 0$. Since f is twice continuously differentiable, one can see that G is continuously differentiable w.r.t. its second argument (u, v) . Consequently G is locally Lipschitz continuous w.r.t. (u, v) and by the Cauchy–Lipschitz Theorem, for each initial condition, there exists a unique local solution to (gVM-DIN-AVD) and thus to (VM-DIN-AVD). We then show that this solution is actually global (in Appendix A.2) by proving the boundedness of (x, y) .

We omit the existence and uniqueness of the solutions of (CN) and (LM) since these are standard results, obtained with similar arguments. We conclude with the following important remark.

Remark 2.2. Thanks to the first-order reformulation (gVM-DIN-AVD), Theorem 2.1 does not need the Lipschitz continuity of $\nabla^2 f$ and not even of ∇f to obtain global existence and uniqueness. In contrast, these smoothness assumptions would have been required if the standard first-order formulation of (VM-DIN-AVD) in phase-space (position-velocity) was used. In fact, (gVM-DIN-AVD) allows for a natural extension of (VM-DIN-AVD) to non-smooth convex functions as we detail in Appendix A.3, and which is not the case for (CN) and (LM). This allows defining an inertial Newton-like dynamics in the non-smooth setting which is important for many applications; e.g. [24].

3. Non-asymptotic Control of (VM-DIN-AVD). The purpose of this section is to understand how close x might be to x_N and x_{LM} , as a function of α and ε . Since f is coercive and strongly convex on bounded sets, it has a unique minimizer $x^* \in \mathbb{R}^n$. Consequently, any two trajectories that converge to x^* will eventually be arbitrarily close to each other. Thus, asymptotic results of the form $\|x(t) - x_N(t)\| \xrightarrow[t \rightarrow +\infty]{} 0$

are not precise enough to claim, for example, that x has a “Newtonian behavior”. Instead, we will derive upper bounds on the distance between trajectories that hold for all time $t \geq 0$, and which typically depend on ε and/or α . We first present the case where α is small relative to ε and then generalize.

3.1. Comparison with (CN) under Moderate Viscous Dampings. When the damping α remains moderate w.r.t. the variable mass ε , one expects the solutions of (VM-DIN-AVD) to be close to that of (CN). We make the following assumptions.

ASSUMPTION 2. *Assumption 1 holds and there exists $c_1, c_2 \geq 0$ such that for all $t \geq 0$, $|\varepsilon'(t)| \leq c_1\varepsilon(t)$ and $\alpha(t) \leq c_2\varepsilon(t)$.*

The assumption states that α must decrease at least as fast as ε (up to a constant).³ The reason for assuming $|\varepsilon'(t)| \leq c_1\varepsilon(t)$ is technical and will appear more clearly in the proofs below. It formally means that ε can decrease at most exponentially fast.⁴ This is a relatively mild assumption that holds, for example, for any polynomial decay $\varepsilon_0/(t+1)^a$, $a \in \mathbb{N}$. We start with the main result of this section.

THEOREM 3.1. *Let x_N be the solution of (CN), and let $c_1, c_2 \geq 0$. There exist $C_0, C_1, C_2 \geq 0$, depending on c_1, c_2 , such that for all (ε, α) for which Assumption 2 holds with constants c_1 and c_2 , the corresponding solution x of (VM-DIN-AVD) is such that for all $t \geq 0$,*

$$(3.1) \quad \|x(t) - x_N(t)\| \leq C_0 e^{-\frac{t}{\beta}} \varepsilon_0 \|\dot{x}_0\| + C_1 \sqrt{\varepsilon(t)} + C_2 \int_{s=0}^t e^{\frac{1}{\beta}(s-t)} \sqrt{\varepsilon(s)} \, ds.$$

This extends a previous result from [2, Proposition 3.1] which states a similar bound for constant ε , $\alpha \equiv 0$ and $\beta = 1$. Theorem 3.1 corresponds to the blue parts in the phase diagram of Figure 1 (see also Corollary 3.6 below).

Remark 3.2. The strength of the result comes from the fact that C_0, C_1, C_2 do not depend on ε and α , and that the result is *non-asymptotic*. This allows in particular for choosing (ε, α) to control the distance from x to x_N , for all time $t \geq 0$.

Remark 3.3. Under Assumption 2, the dynamics (VM-DIN-AVD) is dominated by the variable mass. The damping α does not appear in Theorem 3.1.

The above theorem and remarks emphasize the “Newtonian nature” of (VM-DIN-AVD). We present two lemmas before proving Theorem 3.1, and then state a simpler bound than (3.1), see Corollary 3.6.

LEMMA 3.4. *Let (ε, α) , and let x be the corresponding solution of (VM-DIN-AVD). For all $t \geq 0$, define the function, $U(t) = \frac{\varepsilon(t)}{2} \|\dot{x}(t)\|^2 + f(x(t)) - f(x^*)$. Then U is differentiable and for all $t > 0$,*

$$\frac{dU}{dt}(t) = \frac{\varepsilon'(t)}{2} \|\dot{x}(t)\|^2 - \alpha(t) \|\dot{x}(t)\|^2 - \beta \langle \nabla^2 f(x(t)) \dot{x}(t), \dot{x}(t) \rangle \leq 0.$$

Therefore, in particular, U is non-increasing.

Proof. Let $t \geq 0$, since x is twice differentiable, U is differentiable and,

$$\frac{dU}{dt}(t) = \frac{\varepsilon'(t)}{2} \|\dot{x}(t)\|^2 + \varepsilon(t) \langle \dot{x}(t), \ddot{x}(t) \rangle + \langle \dot{x}(t), \nabla f(x(t)) \rangle.$$

³Assumption 2 can hold only after some $t_0 \geq 0$, we take $t_0 = 0$ for the sake of simplicity.

⁴This is a consequence of Gronwall’s lemma, see e.g., [27].

We use the fact that x is solution of (VM-DIN-AVD), to substitute $\varepsilon(t)\dot{x}(t)$ by its expression, $\frac{dU}{dt}(t) = \frac{\varepsilon'(t)}{2}\|\dot{x}(t)\|^2 - \alpha(t)\|\dot{x}(t)\|^2 - \beta\langle\nabla^2 f(x(t))\dot{x}(t), \dot{x}(t)\rangle$. By assumption ε is non-increasing so for all $t > 0$, $\varepsilon'(t) \leq 0$. Furthermore f is convex so $\langle\nabla^2 f(x(t))\dot{x}(t), \dot{x}(t)\rangle \geq 0$. Hence U is non-increasing. \square

We then state the following bound.

LEMMA 3.5. *For any (ε, α) , the corresponding solution x of (VM-DIN-AVD) is such that for all $t \geq 0$,*

$$\varepsilon(t)\|\dot{x}(t)\| \leq \sqrt{2U(0)}\sqrt{\varepsilon(t)}.$$

From Lemma 3.4, U is non-increasing so $\forall t \geq 0$, $U(t) \leq U(0)$. Then in particular $\frac{\varepsilon(t)}{2}\|\dot{x}(t)\|^2 \leq U(0)$ and the proof follows by multiplying both sides by $\varepsilon(t)$.

Proof of Theorem 3.1. Let (ε, α) as defined in Sections 1.1 and 2, and let x be the corresponding solution of (VM-DIN-AVD). Then, according to Lemma 3.4, for all $t \geq 0$, $U(t) \leq U(0)$, so in particular $f(x(t)) \leq f(x_0) + \frac{\varepsilon_0}{2}\|\dot{x}_0\|^2$. Denoting $M_0 = f(x_0) + \frac{\varepsilon_0}{2}\|\dot{x}_0\|^2$, the set $K_0 = \{y \in \mathbb{R}^n \mid f(y) \leq M_0\}$ is bounded, since f is coercive ($\lim_{\|y\| \rightarrow +\infty} f(y) = +\infty$). So for all $t \geq 0$, $x(t) \in K_0$. Since M_0 (and hence K_0) depends only on ε_0 , x_0 and \dot{x}_0 , we have proved that for any choice (ε, α) , the corresponding solution x of (VM-DIN-AVD) is inside K_0 at all times. Let x_N be the solution of (CN). One can similarly see that for all $t \geq 0$, $f(x_N(t)) \leq f(x_N(0)) = f(x_0) \leq M_0$. So we also have $x_N(t) \in K_0$ for all $t \geq 0$.

Now, fix $c_1, c_2 > 0$, and let (ε, α) such that Assumption 2 is satisfied with constants c_1, c_2 . Let x be the corresponding solution of (VM-DIN-AVD). Since f is strongly convex on bounded sets, it is strongly convex on K_0 . We denote $\mu_{K_0} > 0$ the strong-convexity parameter of f on K_0 . Equivalently, we have that ∇f is strongly monotone on K_0 , that is, $\forall y_1, y_2 \in K_0$,

$$(3.2) \quad \langle \nabla f(y_1) - \nabla f(y_2), y_1 - y_2 \rangle \geq \mu_{K_0} \|y_1 - y_2\|^2.$$

Let $t \geq 0$, since $x(t) \in K_0$ and $x_N(t) \in K_0$, by combining (3.2) with the Cauchy-Schwarz inequality, we deduce that

$$(3.3) \quad \|x(t) - x_N(t)\| \leq \frac{1}{\mu_{K_0}} \|\nabla f(x(t)) - \nabla f(x_N(t))\|.$$

Therefore, it is sufficient to bound the difference of gradients in order to bound $\|x(t) - x_N(t)\|$. First, remark that (CN) can be rewritten as follows: $\frac{d}{dt}\nabla f(x_N(t)) + \frac{1}{\beta}\nabla f(x_N(t)) = 0$. So we can integrate, for all $t \geq 0$,

$$(3.4) \quad \nabla f(x_N(t)) = e^{-\frac{t}{\beta}} \nabla f(x_0).$$

We now turn our attention to $\nabla f(x(t))$, for which we cannot find a closed-form solution in general. We rewrite (VM-DIN-AVD) in the following equivalent form

$$\frac{d}{dt} [\varepsilon(t)\dot{x}(t) + \beta\nabla f(x(t))] + \frac{1}{\beta}\varepsilon(t)\dot{x}(t) + \nabla f(x(t)) = \left(\frac{1}{\beta}\varepsilon(t) + \varepsilon'(t) - \alpha(t) \right) \dot{x}(t).$$

Introducing the variable $\omega(t) = \varepsilon(t)\dot{x}(t) + \beta\nabla f(x(t))$, the latter is thus solution to

$$\dot{\omega}(t) + \frac{1}{\beta}\omega(t) = \left(\frac{1}{\beta}\varepsilon(t) + \varepsilon'(t) - \alpha(t) \right) \dot{x}(t), \quad t \geq 0,$$

with $\omega(0) = \varepsilon_0 \dot{x}_0 + \beta \nabla f(x_0)$. This is a non-homogeneous first-order ODE in ω , whose solution can be expressed using the integrating factor

$$\omega(t) = e^{-\frac{t}{\beta}} (\varepsilon_0 \dot{x}_0 + \beta \nabla f(x_0)) + e^{-\frac{t}{\beta}} \int_0^t e^{\frac{s}{\beta}} \left(\frac{1}{\beta} \varepsilon(s) + \varepsilon'(s) - \alpha(s) \right) \dot{x}(s) ds.$$

We thus have the following expression for $\nabla f(x)$, for all $t \geq 0$,

$$(3.5) \quad \beta \nabla f(x(t)) = \beta e^{-\frac{t}{\beta}} \nabla f(x_0) + e^{-\frac{t}{\beta}} \varepsilon_0 \dot{x}_0 - \varepsilon(t) \dot{x}(t) + e^{-\frac{t}{\beta}} \int_0^t e^{\frac{s}{\beta}} \left(\frac{1}{\beta} \varepsilon(s) + \varepsilon'(s) - \alpha(s) \right) \dot{x}(s) ds.$$

We can now use (3.4) and (3.5) in (3.3) to get

$$\|x(t) - x_N(t)\| \leq \frac{1}{\beta \mu_{\kappa_0}} \left\| e^{-\frac{t}{\beta}} \varepsilon_0 \dot{x}_0 - \varepsilon(t) \dot{x}(t) + e^{-\frac{t}{\beta}} \int_0^t e^{\frac{s}{\beta}} \left(\frac{1}{\beta} \varepsilon(s) + \varepsilon'(s) - \alpha(s) \right) \dot{x}(s) ds \right\|.$$

Using the triangle inequality, we obtain,

$$(3.6) \quad \|x(t) - x_N(t)\| \leq \frac{\varepsilon_0 \|\dot{x}_0\|}{\beta \mu_{\kappa_0}} e^{-\frac{t}{\beta}} + \frac{\varepsilon(t) \|\dot{x}(t)\|}{\beta \mu_{\kappa_0}} + \frac{1}{\beta \mu_{\kappa_0}} \int_0^t e^{\frac{1}{\beta}(s-t)} \left| \frac{1}{\beta} \varepsilon(s) + \varepsilon'(s) - \alpha(s) \right| \|\dot{x}(s)\| ds.$$

The first term in (3.6) corresponds to the first one in (3.1) with $C_0 = 1/(\beta \mu_{\kappa_0})$. As for the second-one, by direct application of Lemma 3.5, for all $t \geq 0$, $\varepsilon(t) \|\dot{x}(t)\| \leq \sqrt{2U(0)} \sqrt{\varepsilon(t)}$, so we set $C_1 = \sqrt{2U(0)}/(\beta \mu_{\kappa_0})$. Regarding the last term in (3.6), using Assumption 2 and again Lemma 3.5, it holds that, for all $s \geq 0$,

$$\left| \frac{1}{\beta} \varepsilon(s) + \varepsilon'(s) - \alpha(s) \right| \|\dot{x}(s)\| \leq \left(\frac{1}{\beta} + c_1 + c_2 \right) \varepsilon(s) \|\dot{x}(s)\| \leq \left(\frac{1}{\beta} + c_1 + c_2 \right) \sqrt{2U(0)} \sqrt{\varepsilon(s)}.$$

This proves the theorem with $C_2 = \frac{\sqrt{2U(0)}}{\beta \mu_{\kappa_0}} \left(\frac{1}{\beta} + c_1 + c_2 \right)$. \square

Let us analyze the bound in Theorem 3.1. The first term in (3.1) decays exponentially fast and can even be zero if the initial speed is $\dot{x}_0 = 0$, the second-one decays like $\sqrt{\varepsilon(t)}$, however, the rate at which the last term decreases is less obvious. The following corollary gives a less-tight but easier-to-understand estimate.

COROLLARY 3.6. *Let the same assumptions and variables as in Theorem 3.1. If furthermore $c_1 < \frac{2}{\beta}$, then there exists $C_3 > 0$ such that for all $t \geq 0$,*

$$\|x(t) - x_N(t)\| \leq C_0 e^{-\frac{t}{\beta}} \varepsilon_0 \|\dot{x}_0\| + C_3 \sqrt{\varepsilon(t)}.$$

Proof of Corollary 3.6. For all $t \geq 0$, define $J(t) = \int_0^t e^{\frac{s}{\beta}} \sqrt{\varepsilon(s)} ds$. An integration by parts yields

$$(3.7) \quad J(t) = \left[\beta e^{\frac{s}{\beta}} \sqrt{\varepsilon(s)} \right]_{s=0}^t - \int_{s=0}^t \beta e^{\frac{s}{\beta}} \frac{\varepsilon'(s)}{2\sqrt{\varepsilon(s)}} ds = \beta e^{\frac{t}{\beta}} \sqrt{\varepsilon(t)} - \beta \varepsilon_0 + \int_{s=0}^t \beta e^{\frac{s}{\beta}} \frac{-\varepsilon'(s)}{2\varepsilon(s)} \sqrt{\varepsilon(s)} ds.$$

By assumption, $0 \leq c_1 < \frac{2}{\beta}$ such that for all $s > 0$, $|\varepsilon'(s)| \leq c_1\varepsilon(s)$, which in our setting is equivalent to $\frac{-\varepsilon'(s)}{\varepsilon(s)} \leq c_1$. So we deduce from (3.7) that

$$J(t) \leq \beta e^{\frac{t}{\beta}} \sqrt{\varepsilon(t)} + c_1 \frac{\beta}{2} \int_{s=0}^t e^{\frac{s}{\beta}} \sqrt{\varepsilon(s)} ds = \beta e^{\frac{t}{\beta}} \sqrt{\varepsilon(t)} + c_1 \frac{\beta}{2} J(t).$$

So, $\left(1 - c_1 \frac{\beta}{2}\right) J(t) \leq \beta e^{\frac{t}{\beta}} \sqrt{\varepsilon(t)}$. By assumption $1 - c_1 \frac{\beta}{2} > 0$, therefore, $J(t) \leq \frac{2}{2 - c_1 \beta} e^{\frac{t}{\beta}} \sqrt{\varepsilon(t)}$. We use this in (3.1) and set $C_3 = C_1 + C_2 \frac{2}{2 - c_1 \beta}$ to get the result. \square

Remark 3.7. Observe that local strong convexity is only required to get (3.2) and (3.3). In fact, local strong convexity can be greatly weakened and our claims above can be generalized to a large sub-class of strictly convex functions as we explain in Appendix B. We did not directly consider the strictly convex case to emphasize the main ideas and ease the reading. Following Remark 2.2, it seems also that a non-smooth extension is possible using regularization techniques. Unlike the strictly convex setting, the non-smooth one would however require further investigations since the vanilla Newton method is not applicable to non-smooth functions.

So far our results only cover the case where α is “not too large” w.r.t. ε , and do not study (LM). We now state a more general result that covers these cases.

3.2. Generalization to Sub-exponentially Decaying Viscous Dampings.

This time we do not assume a link between ε and α but only sub-exponential decays.

ASSUMPTION 3. *Assumption 1 holds and there exists $c_1, c_2 \geq 0$ such that for all $t \geq 0$, $|\varepsilon'(t)| \leq c_1\varepsilon(t)$ and $|\alpha'(t)| \leq c_2\alpha(t)$.*

We are now in position to state the main result of this section.

THEOREM 3.8. *Let x_N and x_{LM} be the solutions of (CN) and (LM) respectively, and let $c_1, c_2 \geq 0$. There exist constants $C, \tilde{C} \geq 0$, depending on $c_1, c_2, \varepsilon_0, \alpha_0$ and the initial conditions, such that for all ε and α for which Assumption 3 holds with c_1, c_2 , the corresponding solution x of (VM-DIN-AVD) is such that for all $t \geq 0$,*

$$(3.8) \quad \|x(t) - x_N(t)\| \leq C \left[e^{-\frac{t}{\beta}} + \sqrt{\varepsilon(t)} + \alpha(t) + \int_{s=0}^t e^{\frac{1}{\beta}(s-t)} (\sqrt{\varepsilon(s)} + \alpha(s)) ds \right],$$

and,

$$(3.9) \quad \|x(t) - x_{LM}(t)\| \leq \tilde{C} \left[e^{-\frac{t}{\beta}} + \sqrt{\varepsilon(t)} + \alpha(t) + \int_{s=0}^t e^{\frac{1}{\beta}(s-t)} (\sqrt{\varepsilon(s)} + \alpha(s)) ds \right].$$

The proof is omitted but key elements are presented in Appendix C. Although it follows a similar reasoning as that of Theorem 3.1, more involved estimates are needed.

Let us comment on these results. The bound (3.8) generalizes Theorem 3.1, although the constant involved will, in general, be larger than those in (3.1) (see the proof of Theorem 3.8 in Appendix C). Theorem 3.8 allows for far more flexibility in the choice of ε and α in order to control x and make it possibly close to x_N . The bound in (3.9) is the same as that in (3.8) (up to a constant), but this time w.r.t. x_{LM} , thus connecting (VM-DIN-AVD) to (LM). We make the following two important remarks. First (3.9) involves α , suggesting that making x close to x_{LM} requires not only ε but also α to vanish asymptotically. Additionally, Theorem 3.8 does not state to which of x_N and x_{LM} the solution of (VM-DIN-AVD) is the closest. It remains an open question to know whether one can make (3.9) independent of α , and to state to which

trajectory x is the closest. Yet, the numerical experiments in Section 5 suggest that neither are possible. Indeed, we observe that for some functions f , x is *sometimes* closer to x_N than to x_{LM} , even when $\varepsilon(t) \leq \alpha(t)$.

Nevertheless, Theorem 3.8 answers the question asked in the introduction: yes, (VM-DIN-AVD) is really of second-order nature since it can be brought close to the second-order dynamics (CN) and (LM). Doing so, it benefits from the good properties of these methods, such as the robustness to bad conditioning, as previously illustrated on the right of Figure 1. This concludes the analysis from a control perspective. We will now derive an approximation of the solution x in order to study the impact that ε and α have on the speed of convergence of x to x^* compared to the speeds of convergence of x_N and x_{LM} .

4. Approximate Solutions and Asymptotic Analysis on Quadratics. In addition to Assumption 1, we consider the case where f is a strongly convex quadratic function in order to study the asymptotic behavior of (VM-DIN-AVD) w.r.t. (CN) and (LM). Quadratic functions are the prototypical example of strongly convex functions. In particular, any strongly convex function can be locally approximated by a quadratic one around its minimizer, making the latter a good model for understanding the local behavior of dynamics. In this section, f is quadratic: $f(y) = \frac{1}{2}\|Ay - b\|_2^2$ for all $y \in \mathbb{R}^n$, where $A \in \mathbb{R}^{n \times n}$ is symmetric positive definite and $b \in \mathbb{R}^n$. Without loss of generality, we take $b = 0$, so that the unique minimum is $x^* = 0$.

4.1. Setting: the Special Case of Quadratic Functions. Quadratic functions are particularly interesting in our setting since DIN-like ODEs take a simpler form (as observed in [8, 46]). Indeed, $\forall y \in \mathbb{R}^n$, $\nabla f(y) = A^T Ay$ and $\nabla^2 f(y) = A^T A$. Since $\nabla^2 f(y)$ is independent of y we can rewrite (VM-DIN-AVD) in an eigenspace⁵ of $A^T A$. That is, we can study the ODE coordinate-wise by looking at one-dimensional problems of the form

$$(Q1\text{-VM-DIN-AVD}) \quad \varepsilon(t)\ddot{x}(t) + (\alpha(t) + \beta\lambda)\dot{x}(t) + \lambda x(t) = 0, \quad t \geq 0.$$

Here (and throughout what follows) $\lambda > 0$ denotes any eigenvalue of $A^T A$ and $x: \mathbb{R}_+ \rightarrow \mathbb{R}$ now denotes the corresponding coordinate (function) of the solution of (VM-DIN-AVD) in an eigenspace of $A^T A$. The dynamics (Q1-VM-DIN-AVD) is a *linear* second-order ODE in x with non-constant coefficients. Similarly, (LM) can be rewritten coordinate-wise as

$$(Q1\text{-LM}) \quad (\alpha(t) + \beta\lambda)\dot{x}_{LM}(t) + \lambda x_{LM}(t) = 0, \quad t \geq 0,$$

where $x_{LM}: \mathbb{R}_+ \rightarrow \mathbb{R}$, and (CN) becomes

$$(Q1\text{-CN}) \quad \beta\dot{x}_N(t) + x_N(t) = 0, \quad t \geq 0,$$

where again, $x_N: \mathbb{R}_+ \rightarrow \mathbb{R}$ is one-dimensional. Observe in particular that (CN) and (LM) are now first-order *linear* ODEs, whose solutions have closed forms: $\forall t \geq 0$,

$$(4.1) \quad x_N(t) = x_0 e^{-\frac{t}{\beta}} \quad \text{and} \quad x_{LM}(t) = x_0 \exp\left(-\int_0^t \frac{\lambda}{\alpha(s) + \beta\lambda} ds\right).$$

Since the minimizer is $x^* = 0$, we see that x_N converges exponentially fast to x^* , with a rate independent of λ while the rate of x_{LM} depends on λ and how fast α vanishes.

⁵This can be generalized to the case where $A^T A$ is only semi-definite by considering orthogonal projections on an eigenspace spanned by the positive eigenvalues of $A^T A$.

Unfortunately, except for some special choices of ε and α (see [8]), one cannot solve the second-order linear ODE (Q1-VM-DIN-AVD) in closed form in general. Additionally, it is hopeless to circumvent the difficulty by finding a closed form for $\nabla f(x)$, accordingly to what we did in Section 3, since here $\nabla f(x) = \lambda x$. In order to study the speed of convergence of x despite not having access to a closed form, we will approximate it with a controlled error, via a method that we now present.

4.2. The Liouville–Green Method. In what follows, we rely on the Liouville–Green method [35, 31], a technique for obtaining *non-asymptotic* approximations to solutions of linear second-order ODEs with non-constant coefficients. First, we give the intuition behind the method, following the presentation of [42]. Consider the differential equation

$$(4.2) \quad \ddot{z}(t) - r(t)z(t) = 0, \quad t \geq 0,$$

where r is real-valued, positive, and twice continuously differentiable. Any linear second-order ODE can be reformulated in the form (4.2), see Lemma 4.5 below. Since for all $t \geq 0$, $r(t) \neq 0$, we can use the changes of variables $\tau = \int_0^t \sqrt{r(s)} \, ds$ and $w = r^{1/4}z$ and show that w is solution to

$$(4.3) \quad \ddot{w}(\tau) - (1 + \psi(\tau))w(\tau) = 0, \quad t \geq 0,$$

where⁶ $\psi(\tau) = \frac{4r(t)r''(t) - 5r'(t)^2}{16r(t)^3}$. The LG method consists in neglecting the term $\psi(\tau)$ in (4.3), which simply yields two approximate solutions $\hat{w}_1(\tau) = e^\tau$ and $\hat{w}_2(\tau) = e^{-\tau}$. Expressing this in terms of z and t , we obtain

$$(4.4) \quad \hat{z}_1(t) = r(t)^{-1/4} \exp\left(\int_0^t \sqrt{r(s)} \, ds\right) \quad \text{and} \quad \hat{z}_2(t) = r(t)^{-1/4} \exp\left(\int_0^t -\sqrt{r(s)} \, ds\right).$$

Those are the LG approximations of the solutions of (4.2). They are formally valid on any $[0, T]$, $T > 0$ when ψ is “not too large” and if \sqrt{r} is integrable on $[0, T]$.

Remark 4.1. There exists other (but less intuitive) ways to derive the LG approximations which allow for generalization to higher-order linear ODEs [17, Chapter 10].

The advantage of this approach is the possibility to estimate the error made using (4.4) w.r.t. the true solutions of (4.2). This is expressed in the following theorem which gathers results from [18, 41, 48].

THEOREM 4.2 ([42]). *Let $r: \mathbb{R}_+ \rightarrow \mathbb{R}$ be a real, positive, twice continuously differentiable function, and define $\varphi(t) = \frac{4r(t)r''(t) - 5r'(t)^2}{16r(t)^{5/2}}$ for all $t \geq 0$. Then for any $T > 0$, the differential equation,*

$$(4.5) \quad \ddot{z}(t) - r(t)z(t) = 0, \quad t \in [0, T],$$

has two real and twice continuously differentiable solutions defined $\forall t \in [0, T]$ by,

$$z_1(t) = \frac{1 + \delta_1(t)}{r(t)^{1/4}} \exp\left(\int_0^t \sqrt{r(s)} \, ds\right) \quad \text{and} \quad z_2(t) = \frac{1 + \delta_2(t)}{r(t)^{1/4}} \exp\left(-\int_0^t \sqrt{r(s)} \, ds\right),$$

where $|\delta_1(t)| \leq \exp\left(\frac{1}{2} \int_0^t |\varphi(s)| \, ds\right) - 1$ and $|\delta_2(t)| \leq \exp\left(-\frac{1}{2} \int_t^T |\varphi(s)| \, ds\right) - 1$.

If in addition $\int_0^{+\infty} |\varphi(s)| \, ds < +\infty$, then the results above also hold for $T = +\infty$.

⁶We express $\psi(\tau)$ via t using the one-to-one correspondence between τ and t to ease readability.

Remark 4.3. We make the following remarks regarding the above result.

- Note that z_1 and z_2 in Theorem 4.2 are *exact* solutions to (4.5). The LG approximations \hat{z}_1 and \hat{z}_2 are obtained by neglecting the unknown functions δ_1 and δ_2 in z_1 and z_2 . The theorem gives a *non-asymptotic* bound for the errors $|z_1(t) - \hat{z}_1(t)|$ and $|z_2(t) - \hat{z}_2(t)|$, $t \geq 0$.
- Since we assumed r to be twice continuously differentiable and positive, φ is continuous, so it is integrable except maybe for $t \rightarrow +\infty$.
- For the sake of simplicity, the formulation of Theorem 4.2 slightly differs from that in [42], the original formulation can be recovered by a change of variable.

4.3. Liouville–Green Approximation of (VM-DIN-AVD). We now proceed to make use of the LG method for approximating the solutions of (Q1-VM-DIN-AVD). The reader only interested in the result can jump directly to the Section 4.4. We first make the following assumption.

ASSUMPTION 4. *The functions α and ε are three times continuously differentiable, and ε_0 is such that $\forall t \geq 0$, $\varepsilon_0 < \frac{(\beta\lambda)^2}{2|\alpha'(t)|+4\lambda}$.*

Remark 4.4. The condition on ε_0 in Assumption 4 is only technical, so that r defined below is positive. It can be easily satisfied since $|\alpha'(t)|$ is uniformly bounded. Indeed, α is non-increasing and non-negative (see Section 1.1), from which one can deduce that $\int_0^{+\infty} |\alpha'(s)| ds \leq \alpha_0$.

We now rewrite (Q1-VM-DIN-AVD) in the form (4.5).

LEMMA 4.5. *Suppose that Assumption 4 holds, and let x be the solution of (Q1-VM-DIN-AVD). For all $t \geq 0$, define*

$$(4.6) \quad p(t) = \frac{\alpha(t) + \beta\lambda}{\varepsilon(t)} \quad \text{and} \quad r(t) = \frac{p(t)^2}{4} + \frac{p'(t)}{2} - \frac{\lambda}{\varepsilon(t)}.$$

Then, p and r are twice continuously differentiable, r is positive and the function y defined for all $t \geq 0$ by $y(t) = x(t) \exp\left(\int_0^t \frac{p(s)}{2} ds\right)$ is a solution to

$$(4.7) \quad \ddot{y}(t) - r(t)y(t) = 0, \quad t \geq 0,$$

with initial condition $(y(0), \dot{y}(0)) = (x_0, \dot{x}_0 + \frac{p(0)}{2}x_0)$.

Proof. We first check that for all $t \geq 0$, $r(t)$ is positive. Let $t > 0$,

$$(4.8) \quad r(t) > 0 \iff \frac{(\alpha(t) + \beta\lambda)^2}{4\varepsilon(t)^2} + \frac{\alpha'(t)}{2\varepsilon(t)} - \frac{(\alpha(t) + \beta\lambda)\varepsilon'(t)}{\varepsilon(t)^2} - \frac{\lambda}{\varepsilon(t)} > 0.$$

Since $\varepsilon'(t) \leq 0$ and $\alpha'(t) \leq 0$, one can check that a sufficient condition for (4.8) to hold is,

$$r(t) > 0 \iff \frac{(\alpha(t) + \beta\lambda)^2}{4} > \left(\frac{|\alpha'(t)|}{2} + \lambda\right) \varepsilon(t) \iff \frac{(\beta\lambda)^2}{2|\alpha'(t)| + 4\lambda} > \varepsilon_0.$$

So under Assumption 4, for all $t \geq 0$, $r(t) > 0$. We then check that y is indeed solution to (4.7). Let $t > 0$,

$$\dot{y}(t) = \frac{p(t)}{2}x(t) \exp\left(\int_0^t \frac{p(s)}{2} ds\right) + \dot{x}(t) \exp\left(\int_0^t \frac{p(s)}{2} ds\right), \quad \text{and,}$$

$$\ddot{y}(t) = \exp\left(\int_0^t \frac{p(s)}{2} ds\right) \left[\left(\frac{p(t)^2}{4} + \frac{p'(t)}{2}\right) x(t) + p(t)\dot{x}(t) + \ddot{x}(t) \right].$$

Since x solves (Q1-VM-DIN-AVD), it holds that $\ddot{x}(t) = -p(t)\dot{x}(t) - \frac{\lambda}{\varepsilon(t)}x(t)$, so,

$$\begin{aligned} \ddot{y}(t) &= \exp\left(\int_0^t \frac{p(s)}{2} ds\right) \left(\frac{p(t)^2}{4} + \frac{p'(t)}{2} - \frac{\lambda}{\varepsilon(t)}\right) x(t) \\ &= \left(\frac{p(t)^2}{4} + \frac{p'(t)}{2} - \frac{\lambda}{\varepsilon(t)}\right) y(t) = r(t)y(t). \quad \square \end{aligned}$$

Lemma 4.5 gives a reformulation of (Q1-VM-DIN-AVD) suited to apply Theorem 4.2. To use the theorem for all $t \geq 0$, we need to ensure that $\varphi(t) = \frac{4r(t)r''(t) - 5r'(t)^2}{16r(t)^{5/2}}$ is integrable. To this aim we make the following assumption.

ASSUMPTION 5. *The functions ε and α have first, second and third-order derivatives that are integrable on $[0, +\infty[$. In addition, $\lim_{t \rightarrow \infty} \varepsilon(t) = 0$ and $\varepsilon'(t)^2/\varepsilon(t)$ is integrable on $[0, +\infty[$.*

Remark 4.6. Assumption 5 holds for most decays used in practice, with in particular any polynomial decay of the form $\frac{\varepsilon_0}{(t+1)^a}$ and $\frac{\alpha_0}{(t+1)^b}$, $a \in \mathbb{N} \setminus \{0\}$ and $b \in \mathbb{N}$. Note that ε and α need not be integrable and α can even be constant.

The next lemma states the integrability of φ on $[0, +\infty[$.

LEMMA 4.7. *Under Assumption 4 and 5, $\int_0^{+\infty} |\varphi(s)| ds < +\infty$.*

The proof of this lemma involves long computations and is thus postponed to Appendix D. We can now use Theorem 4.2 to obtain an exact form for the solution of (Q1-VM-DIN-AVD) based on the LG approximations.

THEOREM 4.8. *Suppose that Assumptions 4 and 5 hold. There exists $A, B \in \mathbb{R}$ such that $x(0) = x_0$, $\dot{x}(0) = \dot{x}_0$ and for all $t \geq 0$, the solution of (Q1-VM-DIN-AVD) is*

$$(4.9) \quad \begin{aligned} x(t) &= A \frac{1 + \delta_1(t)}{r(t)^{1/4}} \frac{\sqrt{\alpha(t) + \beta\lambda}}{\sqrt{\varepsilon(t)}} \exp\left(\int_0^t -\frac{\lambda}{\alpha(s) + \beta\lambda} - \frac{\lambda^2 \varepsilon(s)}{(\alpha(s) + \beta\lambda)^3} + o(\varepsilon(s)) ds\right) \\ &+ B \frac{1 + \delta_2(t)}{r(t)^{1/4}} \frac{\sqrt{\varepsilon(t)}}{\sqrt{\alpha(t) + \beta\lambda}} \\ &\exp\left(\int_0^t -\frac{\alpha(s) + \beta\lambda}{\varepsilon(s)} + \frac{\lambda}{\alpha(s) + \beta\lambda} + \frac{\lambda^2 \varepsilon(s)}{(\alpha(s) + \beta\lambda)^3} + o(\varepsilon(s)) ds\right), \end{aligned}$$

where for all $t \geq 0$,

$$(4.10) \quad |\delta_1(t)| \leq \exp\left(\frac{1}{2} \int_0^t |\varphi(s)| ds\right) - 1 \quad \text{and} \quad |\delta_2(t)| \leq \exp\left(-\frac{1}{2} \int_t^{+\infty} |\varphi(s)| ds\right) - 1.$$

Thanks to the bounds (4.10), we now have an approximation of x . We will use it in particular to compare x asymptotically to the solutions of (Q1-LM) and (Q1-CN). Before this, we prove Theorem 4.8.

Proof of Theorem 4.8. Let x be the solution of (Q1-VM-DIN-AVD) define p, r as in (4.6). Let us also define $y(t) \stackrel{\text{def}}{=} x(t) \exp\left(\int_0^t \frac{p(s)}{2} ds\right)$. According to Lemma 4.5, r is positive and y is solution to (4.7). Then, from Lemma 4.7, $\int_t^T |\varphi(s)| ds < +\infty$, so we can apply Theorem 4.2 to y on $[0, +\infty[$. Therefore, there exists $A, B \in \mathbb{R}$, such that $\forall t \geq 0$,

$$y(t) = A \frac{1 + \delta_1(t)}{r(t)^{1/4}} \exp\left(\int_0^t \sqrt{r(s)} ds\right) + B \frac{1 + \delta_2(t)}{r(t)^{1/4}} \exp\left(\int_0^t -\sqrt{r(s)} ds\right),$$

where A and B are determined by the initial conditions, and δ_1, δ_2 are such that (4.10) holds. Going back to $x(t) = y(t) \exp\left(\int_0^t -\frac{p(s)}{2} ds\right)$, we obtain that for all $t \geq 0$,

$$(4.11) \quad x(t) = A \frac{1 + \delta_1(t)}{r(t)^{1/4}} \exp\left(\int_0^t -\frac{p(s)}{2} + \sqrt{r(s)} ds\right) + B \frac{1 + \delta_2(t)}{r(t)^{1/4}} \exp\left(\int_0^t -\frac{p(s)}{2} - \sqrt{r(s)} ds\right).$$

It now remains to expand the terms in the two exponentials in (4.11) in order to obtain (4.9). To this aim, we approximate $\sqrt{r(s)}$, let $s \geq 0$,

(4.12)

$$\begin{aligned} \sqrt{r(s)} &= \frac{p(s)}{2} \sqrt{1 + \frac{2p'(s)}{p(s)^2} - \frac{4\lambda}{\varepsilon(s)p(s)^2}} \\ &= \frac{p(s)}{2} \left(1 + \frac{p'(s)}{p(s)^2} - \frac{2\lambda}{\varepsilon(s)p(s)^2} - \frac{1}{8} \left(\frac{2p'(s)}{p(s)^2} - \frac{4\lambda}{\varepsilon(s)p(s)^2}\right)^2 + o(\varepsilon(s)^2)\right) \\ &= \frac{p(s)}{2} + \frac{p'(s)}{2p(s)} - \frac{\lambda}{\varepsilon(s)p(s)} - \frac{1}{16} \left(\frac{2p'(s)}{p(s)^{3/2}} - \frac{4\lambda}{\varepsilon(s)p(s)^{3/2}}\right)^2 + o(\varepsilon(s)) \\ &= \frac{p(s)}{2} + \frac{p'(s)\varepsilon(s)}{2(\alpha(s) + \beta\lambda)} - \frac{\lambda}{\alpha(s) + \beta\lambda} - \frac{1}{16} \left(\frac{2p'(s)}{p(s)^{3/2}} - \frac{4\lambda\sqrt{\varepsilon(s)}}{(\alpha(s) + \beta\lambda)^{3/2}}\right)^2 + o(\varepsilon(s)) \\ &= \frac{p(s)}{2} + \frac{\alpha'(s)/2 - \lambda}{\alpha(s) + \beta\lambda} - \frac{\varepsilon'(s)}{2\varepsilon(s)} - \frac{1}{16} \left(\frac{2p'(s)}{p(s)^{3/2}} - \frac{4\lambda\sqrt{\varepsilon(s)}}{(\alpha(s) + \beta\lambda)^{3/2}}\right)^2 + o(\varepsilon(s)) \end{aligned}$$

To ease the readability, we denote $h(t) = \left(\frac{2p'(s)}{p(s)^{3/2}} - \frac{4\lambda\sqrt{\varepsilon(s)}}{(\alpha(s) + \beta\lambda)^{3/2}}\right)^2$. Focusing on the first exponential term in (4.11), we deduce from (4.12) that for all $t \geq 0$,

$$\begin{aligned} &\exp\left(\int_0^t -\frac{p(s)}{2} + \sqrt{r(s)} ds\right) \\ &= \exp\left(\int_0^t \frac{\alpha'(s)/2 - \lambda}{\alpha(s) + \beta\lambda} - \frac{\varepsilon'(s)}{2\varepsilon(s)} - \frac{\lambda}{\alpha(s) + \beta\lambda} - \frac{1}{16} h(t) + o(\varepsilon(s)) ds\right) \\ &= \frac{\sqrt{\alpha(t) + \beta\lambda}}{\sqrt{\alpha_0 + \beta\lambda}} \frac{\sqrt{\varepsilon_0}}{\sqrt{\varepsilon(t)}} \exp\left(\int_0^t \frac{-\lambda}{\alpha(s) + \beta\lambda} - \frac{1}{16} h(t) + o(\varepsilon(s)) ds\right) \\ &= \frac{\sqrt{\alpha(t) + \beta\lambda}}{\sqrt{\alpha_0 + \beta\lambda}} \frac{\sqrt{\varepsilon_0}}{\sqrt{\varepsilon(t)}} \exp\left(\int_0^t \frac{-\lambda}{\alpha(s) + \beta\lambda} - \frac{\lambda^2 \varepsilon(s)}{(\alpha(s) + \beta\lambda)^3} + o(\varepsilon(s)) ds\right), \end{aligned}$$

where the last line relies on further computations postponed to Lemma D.1 in Appendix D. Performing the exact same type of computations on $\exp\left(\int_0^t -\frac{p(s)}{2} - \sqrt{r(s)} ds\right)$, and up to redefining A and B so as to encompass all the constants, we obtain (4.9) and the result is proved. \square

4.4. Comparison of x with x_{LM} and x_N . We now have an expression for x which is almost explicit: we do not know δ_1 and δ_2 in closed form, but they are uniformly bounded (by Lemma 4.7). We will now compare the asymptotic behavior of (4.9) with those of the solutions of (Q1-LM) and (Q1-CN) that we denoted x_{LM} and x_N respectively. Our main result of Section 4 is the following, where $\sim_{+\infty}$ denotes the asymptotic equivalence⁷ between two functions as $t \rightarrow \infty$.

THEOREM 4.9. *Let x be the solution of (Q1-VM-DIN-AVD), given in (4.9), and x_{LM} and x_N whose closed forms are stated in (4.1). Under Assumptions 4 and 5, there exists $C > 0$ such that the following asymptotic equivalences hold:*

$$(4.13) \quad \begin{aligned} x(t) &\sim_{+\infty} x_{LM}(t)C \exp\left(\int_0^t -\frac{\lambda^2\varepsilon(s)}{(\alpha(s) + \beta\lambda)^3} + o(\varepsilon(s)) ds\right), \quad \text{and} \\ x(t) &\sim_{+\infty} x_N(t)C \exp\left(\int_0^t \frac{\alpha(s)}{\beta(\alpha(s) + \beta\lambda)} - \frac{\lambda^2\varepsilon(s)}{(\alpha(s) + \beta\lambda)^3} + o(\varepsilon(s)) ds\right). \end{aligned}$$

As a consequence, the convergence of x to x^* is:

- (i) Faster than that of x_{LM} if ε is non-integrable and as fast otherwise.
- (ii) Slower than that of x_N if α is non-integrable and as fast if α is integrable, in the case where $\forall t \geq 0, \alpha(t) > \varepsilon(t)$.
- (iii) Faster than that of x_N if ε is non-integrable and as fast if ε is integrable, in the case where $\forall t \geq 0, \alpha(t) < \varepsilon(t)$.

While the results of Section 3 were related to the closeness of (VM-DIN-AVD) w.r.t. (CN) and (LM) from a control perspective, Theorem 4.9 provides a different type of insight. First, the results are asymptotic, so they only allow controlling (VM-DIN-AVD) for large t . They provide however a clear understanding of the nature of the solutions of (VM-DIN-AVD) and their convergence. The conclusions (summarized in Table 1) are in accordance with what we would expect: when the viscous damping is larger than the variable mass, (VM-DIN-AVD) behaves more like the Levenberg–Marquardt method than the Newton one, but it actually becomes an accelerated Levenberg–Marquardt dynamics when ε is non-integrable but vanishing. However, when the variable mass ε is larger than α , the dynamics is closer to the one of the Newton method, and can actually be an accelerated Newton dynamics, again for non-integrable ε . This is analogous to the necessary condition that α must be non-integrable in order to accelerate first-order methods in convex optimization (see [7]). We conclude this section by proving Theorem 4.9.

Proof of Theorem 4.9. Thanks to Assumptions 4 and 5, Theorem 4.8 tells us that x has the form (4.9). We now analyze the two terms in (4.9).

First, we know from Theorem 4.8 that $\delta_1(0) = 0$ and $\lim_{t \rightarrow +\infty} \delta_2(t) = 0$. In addition, by Lemma 4.7, δ_1 and δ_2 are uniformly bounded by some positive constant. Then $r(t)^{-1/4}$ decays asymptotically like $\sqrt{\varepsilon(t)}$ and α is bounded. So $A \frac{1+\delta_1(t)}{r(t)^{1/4}} \frac{\sqrt{\alpha(t)+\beta\lambda}}{\sqrt{\varepsilon(t)}}$ is

⁷Two real-valued functions g_1 and g_2 are asymptotically equivalent in $+\infty$ if and only if $\lim_{t \rightarrow \infty} \frac{g_1(t)}{g_2(t)} = 1$.

asymptotically equivalent to some constant $c_1 \in \mathbb{R}$ as $t \rightarrow +\infty$. Similarly, the factor $B \frac{1+\delta_2(t)}{r(t)^{1/4}} \frac{\sqrt{\varepsilon(t)}}{\sqrt{\alpha(t)+\beta\lambda}}$ is equivalent to $c_2\varepsilon(t)$, with $c_2 \in \mathbb{R}$.

We now analyze the ‘‘exponential factors’’ in (4.9). On the one hand, $\frac{\lambda}{\alpha(s)+\beta\lambda} + \frac{\lambda^2\varepsilon(s)}{(\alpha(s)+\beta\lambda)^3} + o(\varepsilon(s))$ converges to $\frac{1}{\beta}$ as $s \rightarrow \infty$, while $\frac{\alpha(s)+\beta\lambda}{\varepsilon(s)}$ diverges to $+\infty$. Therefore, we deduce that,

$$\begin{aligned} \exp\left(\int_0^t -\frac{\alpha(s)+\beta\lambda}{\varepsilon(s)} + \frac{\lambda}{\alpha(s)+\beta\lambda} + \frac{\lambda^2\varepsilon(s)}{(\alpha(s)+\beta\lambda)^3} + o(\varepsilon(s)) \, ds\right) \\ = o\left(\exp\left(\int_0^t -\frac{\lambda}{\alpha(s)+\beta\lambda} - \frac{\lambda^2\varepsilon(s)}{(\alpha(s)+\beta\lambda)^3} + o(\varepsilon(s)) \, ds\right)\right). \end{aligned}$$

As a consequence, the second term in (4.9) will decrease to 0 faster than the first-one (let alone the additional $\varepsilon(t)$ decay that we have just discussed). The asymptotic behavior of x will thus be governed by the first term in (4.9).

Let us now focus on the first term in (4.9). Observe that $\exp\left(\int_0^t -\frac{\lambda}{\alpha(s)+\beta\lambda} \, ds\right)$ is exactly the decay of x_{LM} in (4.1). Thus, we have proved that there exists $C > 0$, such that the following asymptotic equivalence holds,

$$\begin{aligned} A \frac{1+\delta_1(t)}{r(t)^{1/4}} \frac{\sqrt{\alpha(t)+\beta\lambda}}{\sqrt{\varepsilon(t)}} \exp\left(\int_0^t -\frac{\lambda}{\alpha(s)+\beta\lambda} - \frac{\lambda^2\varepsilon(s)}{(\alpha(s)+\beta\lambda)^3} + o(\varepsilon(s)) \, ds\right) \\ \sim_{+\infty} x_{LM}(t)C \exp\left(\int_0^t -\frac{\lambda^2\varepsilon(s)}{(\alpha(s)+\beta\lambda)^3} + o(\varepsilon(s)) \, ds\right), \end{aligned}$$

which proves the first part of (4.13). The second equivalence in (4.13) is obtained using the following identity,

$$(4.14) \quad \int_0^t -\frac{\lambda}{\alpha(s)+\beta\lambda} \, ds = \int_0^t -\frac{1}{\beta} + \frac{\alpha(s)}{\beta(\alpha(s)+\beta\lambda)} \, ds = -\frac{t}{\beta} + \int_0^t \frac{\alpha(s)}{\beta(\alpha(s)+\beta\lambda)} \, ds$$

and $e^{-t/\beta}$ is precisely the rate at which x_N decreases. So (4.13) holds.

It finally remains to deduce the conclusions of the theorem from (4.13).

– Regarding the comparison with x_{LM} , the integral $\int_0^t -\frac{\lambda^2\varepsilon(s)}{(\alpha(s)+\beta\lambda)^3} + o(\varepsilon(s)) \, ds$ converges if and only if ε is integrable on $[0, +\infty[$, and diverges to $-\infty$ when ε is not. So x converges to 0 at least as fast as x_{LM} and faster when ε is not integrable.

– As for the comparison with x_N , if $\alpha(s) > \varepsilon(s) \geq 0$ for all $s \geq 0$, then the integral $\int_0^t \frac{\alpha(s)}{\beta(\alpha(s)+\beta\lambda)} - \frac{\lambda^2\varepsilon(s)}{(\alpha(s)+\beta\lambda)^3} + o(\varepsilon(s)) \, ds$ is convergent in $+\infty$ if and only if α is integrable and diverges to $+\infty$ when α is non-integrable. So when α is integrable, the speed of convergence of x is the same as that of x_N . When α is not integrable, the convergence to 0 is slower but still holds. Indeed, for all $s \geq 0$ $\frac{\alpha(s)}{\beta(\alpha(s)+\frac{1}{\beta})} < \frac{\alpha(s)}{\beta\alpha(s)} = \frac{1}{\beta}$. Thus for

all $t > 0$, $-\frac{t}{\beta} + \int_0^t \frac{\alpha(s)}{\beta(\alpha(s)+\beta\lambda)} \, ds < 0$.

– Finally, the comparison with x_N in the case $\varepsilon(s) > \alpha(s)$ is exactly the same as the comparison with x_{LM} using (4.14). \square

Remark 4.10. As discussed in the beginning of this section, the main motivation for studying quadratic functions is that strongly convex functions with Lipschitz continuous gradient can be tightly both upper- and lower-bounded by quadratic functions around their minimizers. The intuitive attempt to extend our analysis to this class of

functions consists therefore in applying Theorem 4.8 to the upper and lower quadratic bounds. Such analysis is hence only possible for initial conditions close-enough to the minimizer, and even in that case we did not manage to derive precise-enough asymptotic estimates for the solution of (VM-DIN-AVD) using this strategy. Whether this is possible or not remains an open question.

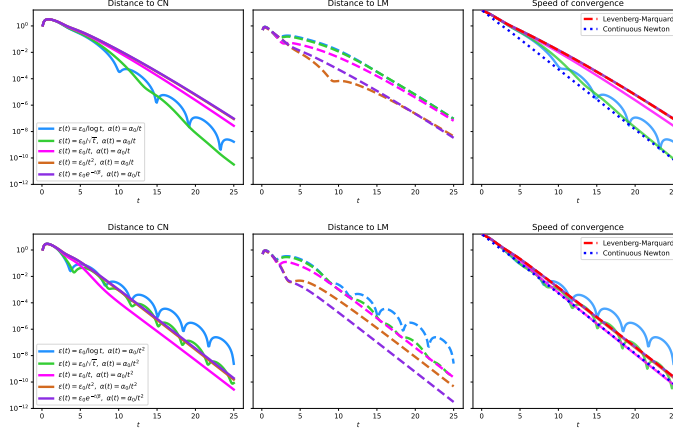


FIG. 2. Comparison of the solutions x_N , x_{LM} and x of (CN), (LM) and (VM-DIN-AVD) respectively, for a strongly convex function of the form $f(x) = e^{-\|x\|^2} + \frac{1}{2}\|Ax\|^2$. Left figures: distance $\|x(t) - x_N(t)\|$ versus time t , each curve corresponds to a different choice of ε ; middle figures: distance $\|x(t) - x_{LM}(t)\|$, again for several ε . Right figures: distance to the optimum x^* for reference, x_N and x_{LM} are in dotted and dashed lines, other curves correspond to (VM-DIN-AVD) for several choices of ε . The brown curve is often hidden behind the purple (and sometimes the pink) curve. Top and bottom rows show results respectively for non-integrable and integrable viscous dampings α . The theoretical bounds from Theorem 3.8 are only displayed on Figure 4 below, for the sake of readability.

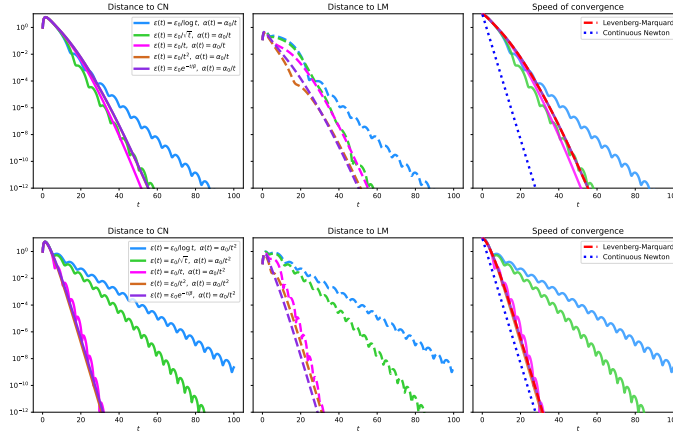


FIG. 3. Similar experiment and figures as those described in Figure 2, but for the function $f(x) = \log(\sum_{i=1}^n e^{x_i} + e^{-x_i}) + \frac{1}{2}\|Ax\|^2$.

5. Numerical Experiments. We present two set of experiments that illustrate our main results from Sections 3 and 4. We first detail the general methodology.

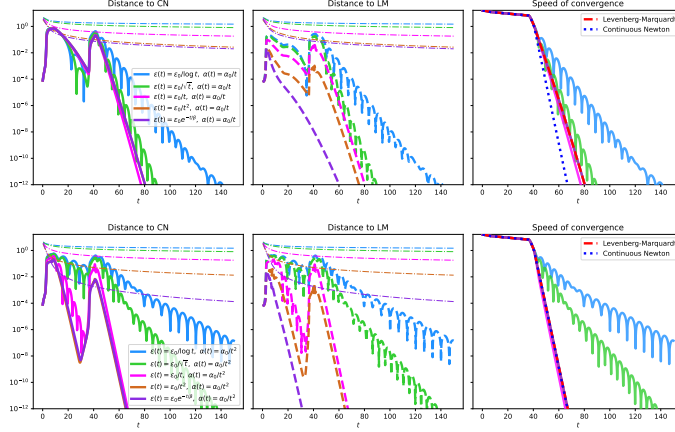


FIG. 4. *Similar experiment and figures as those described in Figure 2, but for the function $f(x) = \sum_{i=1}^n x_i^{50} + \frac{1}{2}\|Ax\|^2$. The thin “dash dotted” curves represent approximations of the theoretical bounds from Theorem 3.8 for each choice of (ϵ, α) considered.*

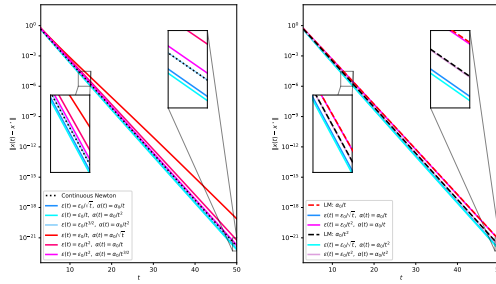


FIG. 5. *Numerical validation of Theorem 4.9: distance to the optimum x^* as a function of time on a quadratic function $f(x) = \frac{1}{2}\|Ax\|^2$. Left: speed comparison w.r.t. (CN) for several choices of ϵ and α . Right: Comparison with LM for α integrable or not and several choices of ϵ . Shades of blue represent cases where $\epsilon(t) > \alpha(t)$ while shades of red represent the opposite setting.*

5.1. Methodology. We compare the solutions of (CN), (LM) and (VM-DIN-AVD) obtained for strongly convex functions in dimension $n = 100$. Since closed-form solutions are not available, they are estimated via discretization schemes with small⁸ step-sizes $\gamma = 10^{-1}$. All ODEs are initialized at $x_0 \in \mathbb{R}^n$, where each coordinate of x_0 belongs to $\{-1, 1\}$. For all functions this means that x_0 is approximately at distance \sqrt{n} of x^* , hence not very close to x^* . We then approximated the ODEs on time intervals that are long-enough to observe both non-asymptotic and asymptotic convergence behaviors. We used Euler semi-explicit schemes obtained by solving linear systems, for the sake of stability. The resulting algorithms are detailed in Appendix E.

5.2. First Experiment: Distance between Trajectories. We begin with an empirical validation of the results of Section 3 on the distance between x , x_{LM} and x_N . Each of Figures 2, 3 and 4 corresponds to a different strongly convex function, specified below its corresponding figure. To ensure strong convexity, each function

⁸In practice Newton’s method would converge faster when used with $\gamma = \beta = 1$, here we take a smaller step-size γ to more accurately approximate the solutions of the ODEs.

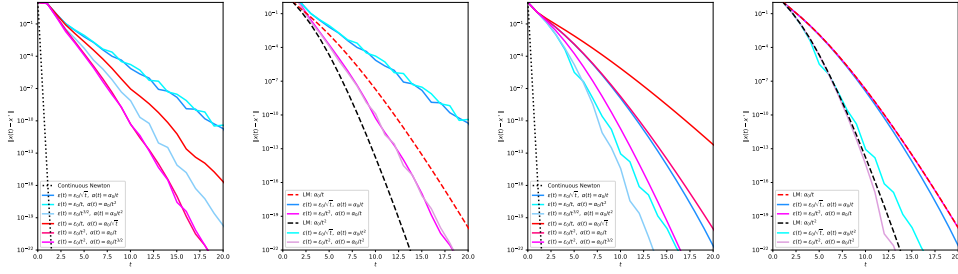


FIG. 6. Same experiments as in Figure 5 but in the setting $\gamma = \beta = 1$ (so that Newton's method converges in one iteration), for large values of ε_0 (left figures) and small ε_0 (right). The use of a large step-size γ makes this setting out of the scope of our theoretical results. Due to fast convergence, different behaviors can be observed for large values of ε_0 and smaller ones.

contains a quadratic term of the form $\|Ax\|^2$, where A is symmetric positive definite.

Several observations can be made from the numerical results, but we first note on the right plots of Figures 2 to 4 that x_N always converges asymptotically linearly (i.e., exponentially fast). This is also the case for x and x_{LM} in some (but not all) cases. This is important because $\|x(t) - x_N(t)\| \leq \|x(t) - x^*\| + \|x_N(t) - x^*\|$, so if both x and x_N converge linearly, then the bounds of Theorems 3.1 and 3.8 need not be asymptotically tight. That being said, the strength of these results is to be non-asymptotic and this is highlighted by the experiments as we now explain.

The left and middle plots of Figures 2, 3 and 4 are consistent with Theorems 3.1 and 3.8, since the distances $\|x(t) - x_N(t)\|$ and $\|x(t) - x_{LM}(t)\|$ decrease relatively fast to zero. Again, when x converges rapidly to x^* this is not very insightful, however, the main interest of our theorems appears on the left of Figures 3 and 4: the blue and green curves, corresponding to slowly decaying choices of ε , converge more slowly than other trajectories. However, when taking faster decays, we recover fast convergence and closeness to x_N (this is particularly true for the purple curve). Very similar observations are made w.r.t. x_{LM} on the middle plots. Despite not being stated in the theorems of Section 3, the experiments match the intuition that when $\varepsilon > \alpha$, x may be closer to x_N and when $\varepsilon \leq \alpha$, x would rather be closer to x_{LM} . This is more noticeable on the top rows of the figures, where α is not integrable.

Figure 4 suggests that the bounds in Theorem 3.8 are rather tight for small t , since, for example, the blue and green curves on the left show a relatively slow vanishing of $\|x(t) - x_N(t)\|$ for slowly decaying ε . The bounds seem however often too pessimistic for large t , for which the second part of our study provides better insights (see Section 4 and below). Interestingly, slow decays of ε might result in faster convergence for x than fast decays (and also faster convergence than x_{LM}), notably on Figure 2. We also note that $\varepsilon(t) = \varepsilon_0/t$ combined either with $\alpha(t) = \alpha_0/t$ or $\alpha(t) = \alpha_0/t^2$ seems to very often yield fast convergence in these experiments.

5.3. Second Experiment: Empirical Validation of Theorem 4.9. We now turn our attention to the solutions x , x_N and x_{LM} for a quadratic function of the form $f(y) = \frac{1}{2}\|Ay\|^2$, $y \in \mathbb{R}^n$, and for several choices of ε and α . The results in Figure 5 exactly match the expected behavior summarized in Table 1. Indeed, looking first at the right-hand side of Figure 5, x is as fast as the corresponding⁹ x_{LM} when ε is not integrable and regardless of α , and x is faster when ε is non-integrable. Then on the

⁹That is, the solution of (LM) for the same α as that considered for (VM-DIN-AVD).

left-hand side, when comparing to x_N , x is slower in settings where α is larger than ε and non-integrable (red curves), or almost as fast when α is integrable (pink curve). However, acceleration w.r.t. to x_N is indeed achieved for non-integrable ε regardless of α (first-two blue curves), and the rate is the same as that of x_N when ε is integrable (third blue curve). Finally, on Figure 6 we empirically investigate the case of large step-sizes γ and take $\gamma = \beta = 1$ (optimal for Newton's method on quadratics). Due to large step-sizes our theoretical results do not hold and the choice of ε_0 matters more. This evidences a trade-off between ε_0 and γ in the case of discrete algorithms which is not present in the analysis of ODEs and is worth investigating in a future work.

6. Conclusions and Perspectives. We introduced a general ODE (VM-DIN-AVD) featuring variable mass, and provided a deep understanding on the behavior of its solutions w.r.t. time dependent control parameters ε and α , both, asymptotically and non-asymptotically. We can conclude that (VM-DIN-AVD) is indeed of (regularized) Newton type, since it can be controlled to be close to both (CN) and (LM). Yet we also showed that (VM-DIN-AVD) fundamentally differs from the other two dynamics in its nature. In particular, Theorem 4.9 and the numerical experiments emphasized that ε and α can accelerate (or slow down) (VM-DIN-AVD) w.r.t. (CN) and (LM). We also note that our bounds in Theorems 3.1 and 3.8 seem relatively tight, in particular for functions with large gradients (see Figure 4). Our contribution yields a complete and satisfying picture on the relation between the three systems, which was only partially understood. We believe that our results build a strong foundation for the development of algorithms that combine the best properties of first- and second-order optimization methods.

As for future work, we showed that (VM-DIN-AVD) is promising from an optimization perspective. So far we approximated solutions of (VM-DIN-AVD) via schemes that required solving linear systems (this is also true for (CN) and (LM)). Our new understanding on (ε, α) paves the way towards designing new Newton-like algorithms with a significantly reduced computational cost, which is crucial for large-scale optimization. Another open question is whether it is possible to preserve the properties evidenced in this work when ε is defined in a closed-loop manner (formally depending on x rather than on t). Finally, it would be worth investigating how the current work can be extended to general convex and/or non-smooth functions.

Acknowledgment. We thank the development teams of the following libraries: Python [44], Numpy [49] and Matplotlib [32].

Appendix A. Equivalent First-order System and Existence of Solutions.

A.1. First-order Equivalent Formulation. We reformulate (VM-DIN-AVD) as a system of ODE involving only first-order time derivatives and the gradient of f . For this purpose, notice that for all $t > 0$ (VM-DIN-AVD) can be rewritten as,

$$(A.1) \quad \frac{d}{dt} [\varepsilon(t)\dot{x}(t)] + \beta \frac{d}{dt} \nabla f(x(t)) + \alpha(t)\dot{x}(t) - \varepsilon'(t)\dot{x}(t) + \nabla f(x(t)) = 0, \quad t \geq 0.$$

We then integrate (A.1) for all $t \geq 0$,

$$(A.2) \quad \varepsilon(t)\dot{x}(t) + \beta \nabla f(x(t)) - \varepsilon_0 \dot{x}_0 - \beta \nabla f(x_0) + \int_0^t (\alpha(s) - \varepsilon'(s))\dot{x}(s) + \nabla f(x(s)) ds = 0.$$

For all $t \geq 0$, we define the variable,

$$z(t) = \int_0^t (\alpha(s) - \varepsilon'(s))\dot{x}(s) + \nabla f(x(s)) \, ds - \varepsilon_0 \dot{x}_0 - \beta \nabla f(x_0).$$

We differentiate z , for all $t > 0$, $\dot{z}(t) = (\alpha(t) - \varepsilon'(t))\dot{x}(t) + \nabla f(x(t))$, so that we can rewrite (A.2) as,

$$\begin{cases} \varepsilon(t)\dot{x}(t) + \beta \nabla f(x(t)) + z(t) = 0 \\ \dot{z}(t) - (\alpha(t) - \varepsilon'(t))\dot{x}(t) - \nabla f(x(t)) = 0 \end{cases}, \quad t \geq 0.$$

We substitute the first line in the second-one,

$$(A.3) \quad \begin{cases} \varepsilon(t)\dot{x}(t) + \beta \nabla f(x(t)) + z(t) = 0 \\ \beta \dot{z}(t) - \beta(\alpha(t) - \varepsilon'(t) - \frac{1}{\beta}\varepsilon(t))\dot{x}(t) + z(t) = 0 \end{cases}, \quad t \geq 0.$$

To ease the readability, we recall the notation $\nu(t) = \alpha(t) - \varepsilon'(t) - \frac{1}{\beta}\varepsilon(t)$ from Section 2. Then define for all $t \geq 0$, $y(t) = z(t) - \nu(t)x(t)$, and differentiate, $\dot{y}(t) = \dot{z}(t) - \nu(t)\dot{x}(t) - \nu'(t)x(t)$. We finally rewrite (A.3) as,

$$\begin{cases} \varepsilon(t)\dot{x}(t) + \beta \nabla f(x(t)) + \nu(t)x(t) + y(t) = 0 \\ \dot{y}(t) + \nu'(t)x(t) + \frac{\nu(t)}{\beta}x(t) + \frac{1}{\beta}y(t) = 0 \end{cases}.$$

which is (gVM-DIN-AVD). Finally, the initial condition on y is

$$y(0) = z(0) - \nu(0)x(0) = -\varepsilon_0 \dot{x}_0 - \beta \nabla f(x_0) - (\alpha_0 - \varepsilon'_0 - \frac{1}{\beta}\varepsilon_0)x_0.$$

Remark A.1. Notice that the quantity $\nu(t) = \alpha(t) - \varepsilon'(t) - \frac{1}{\beta}\varepsilon(t)$ involved in (gVM-DIN-AVD) also plays a key role in our analysis of Section 3, see e.g., (3.6). In particular the sign of $\nu(t)$ changes the nature of (VM-DIN-AVD) and is related to Assumption 2.

A.2. Local Solutions are Global. Using the formulation (gVM-DIN-AVD), we proved local existence and uniqueness of solutions of (VM-DIN-AVD) in Section 2. Using the same notations, we justify that the local solution (x, y) actually exists globally. According to Lemma 3.4, the Lyapunov function $U(t) = \frac{\varepsilon(t)}{2}\|\dot{x}(t)\|^2 + f(x(t)) - f(x^*)$ is non-negative and decreasing. Thus, it is uniformly bounded on \mathbb{R}_+ and the same holds for $t \mapsto f(x(t))$ since for all $t \geq 0$, $U(t) \geq f(x(t))$. Then, f is coercive by assumption, so x is uniformly bounded on \mathbb{R}_+ (otherwise $f(x)$ could not remain bounded). We now prove that y is also uniformly bounded. From (gVM-DIN-AVD), for all $t > 0$, $\dot{y}(t) = -\frac{1}{\beta}y(t) - (\frac{\nu(t)}{\beta} + \nu'(t))x(t)$ so we can use the following integrating factor,

$$y(t) = e^{-\frac{t}{\beta}}y_0 - e^{-\frac{t}{\beta}} \int_0^t \frac{1}{\beta} e^{\frac{s}{\beta}} (\nu(s) + \beta \nu'(s))x(s) \, ds.$$

Using triangle inequalities, for all $t \geq 0$,

$$(A.4) \quad \begin{aligned} \|y(t)\| &\leq e^{-\frac{t}{\beta}}\|y_0\| + \sup_{s \geq 0} \|(\nu(s) + \beta \nu'(s))x(s)\| e^{-\frac{t}{\beta}} \int_0^t \frac{1}{\beta} e^{\frac{s}{\beta}} \, ds \\ &\leq \|y_0\| + \sup_{s \geq 0} \|(\nu(s) + \beta \nu'(s))x(s)\|. \end{aligned}$$

Using the definition of ε and α from Sections 1.1 and 2, observe that ε , α , ε' and α' are all bounded on \mathbb{R}_+ , and ε'' is assumed to be bounded. So ν and ν' are bounded, and since we also proved that x is uniformly bounded on \mathbb{R}_+ , we deduce from (A.4) that y is uniformly bounded as well. Hence, the unique local solution (x, y) is global.

A.3. Existence and Uniqueness for Non-smooth Functions. Following Remark 2.2, we here provide the main arguments to extend (VM-DIN-AVD) to the non-smooth setting. Assume, in this section only, that f is convex, proper, lower semi-continuous, and (without loss of generality) that $\text{dom}(f) = \mathbb{R}^n$. Consider

$$(A.5) \quad \begin{cases} \varepsilon(t)\dot{x}(t) + \beta\partial f(x(t)) + \nu(t)x(t) + y(t) & \ni 0, \\ \dot{y}(t) + \nu'(t)x(t) + \frac{\nu(t)}{\beta}x(t) + \frac{1}{\beta}y(t) & = 0, \end{cases} \quad \text{for a.e. } t \geq 0,$$

where ∂f is the sub-differential of f . Note that when f is differentiable, (A.5) boils down to (gVM-DIN-AVD). An absolutely continuous mapping $(x, y) : \mathbb{R}_+ \rightarrow \mathbb{R}^n \times \mathbb{R}^n$ that satisfies (A.5) is called solution. Remark that defining $\tilde{G}(t, x(t), y(t)) = (\frac{\nu(t)x(t)+y(t)}{\varepsilon(t)}, \nu'(t)x(t) + \frac{\nu(t)}{\beta}x(t) + \frac{1}{\beta}y(t))$ and $F(t, x(t), y(t)) = (\frac{\beta}{\varepsilon(t)}f(x(t)), 0)$, (A.5) rewrites $(\dot{x}(t), \dot{y}(t)) + \partial F(t, x(t), y(t)) + \tilde{G}(t, x(t), y(t)) \ni 0$. Then under our assumptions on ε and α , $\tilde{G}(t, \cdot, \cdot)$ is Lipschitz continuous which allows showing existence of a solution via the theory of non-linear semigroups; see [20, Proposition 3.12]. Note that the chain rule on (x, y) holds for a.e. $t \geq 0$ thanks to absolute continuity, which allows extending our Lyapunov analysis as well.

Appendix B. Generalization to Strictly Convex Functions. Following Remark 3.7, we detail how to generalize Section 3 beyond strongly convex functions. Assume, here only, that f is strictly convex. Since f is continuous and the trajectories are contained in some compact set K_0 (see the proof of Theorem 3.1), it follows that f is uniformly convex on K_0 ; see [15, Proposition 10.15]. This implies that ∇f is uniformly monotone [15, Example 22.3] on K_0 . That is, there exist an increasing function $\phi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ such that $\phi(0) = 0$, and $\forall y_1, y_2 \in K_0$, $\langle \nabla f(y_1) - \nabla f(y_2), y_1 - y_2 \rangle \geq \phi(\|y_1 - y_2\|)$. Note that when f is strongly convex, $\phi(s) = \mu_{K_0}s^2$, $\forall s \geq 0$. We deduce that $\frac{\phi(\|y_1 - y_2\|)}{\|y_1 - y_2\|} \leq \|\nabla f(y_1) - \nabla f(y_2)\|$, which allows extending (3.3) and our analysis to strictly convex f whose associated ϕ is such that $\phi(s)/s \xrightarrow{s \rightarrow 0} 0$. In that case, ϕ is called a growth function, and this covers many applications; see [26, 40].

Appendix C. Proof of Theorem 3.8: Key arguments. This section is devoted to proving the general result of Section 3. Fix some constants $c_1, c_2 > 0$ and let ε and α such that Assumption 3 is satisfied with these constants. Let x be the corresponding solution of (VM-DIN-AVD), x_N , and x_{LM} that of (CN) and (LM), respectively. Following the same arguments as in the beginning of the proof of Theorem 3.1, for all $t \geq 0$, $x(t)$, $x_N(t)$ and $x_{LM}(t)$ belong to the bounded set K_0 defined in that proof. Since f is μ -strongly convex on K_0 , the proof relies again on bounding differences of gradients like in (3.3). Following the exact same steps as in the proof of Theorem 3.1, we know the closed form of $\nabla f(x_N)$ given in (3.4), and an expression for $\nabla f(x)$ given in (3.5). For all $s \geq 0$, we have the identity,

$$(C.1) \quad e^{\frac{s}{\beta}}\dot{x}(s) = e^{\frac{s}{\beta}}\dot{x}(s) + \frac{1}{\beta}e^{\frac{s}{\beta}}x(s) - \frac{1}{\beta}e^{\frac{s}{\beta}}x(s) = \frac{d}{ds}(e^{\frac{s}{\beta}}x(s)) - \frac{1}{\beta}e^{\frac{s}{\beta}}x(s),$$

which we use to perform an integration by parts

$$\int_0^t e^{\frac{s}{\beta}}\alpha(s)\dot{x}(s) ds = \left[\alpha(s)e^{\frac{s}{\beta}}x(s)\right]_0^t - \int_0^t \left(\alpha'(s) + \frac{\alpha(s)}{\beta}\right) e^{\frac{s}{\beta}}x(s) ds.$$

Therefore,

(C.2)

$$e^{-\frac{t}{\beta}} \int_0^t e^{\frac{s}{\beta}} \alpha(s) \dot{x}(s) ds = \alpha(t)x(t) - e^{-\frac{t}{\beta}} \alpha_0 x_0 - \int_0^t e^{\frac{s-t}{\beta}} \left(\alpha'(s) + \frac{\alpha(s)}{\beta} \right) x(s) ds,$$

and we can substitute in (3.5),

(C.3)

$$\begin{aligned} \beta \nabla f(x(t)) &= \beta e^{-\frac{t}{\beta}} \nabla f(x_0) + e^{-\frac{t}{\beta}} \varepsilon_0 \dot{x}_0 - \varepsilon(t) \dot{x}(t) + \int_0^t e^{\frac{s-t}{\beta}} \left(\frac{1}{\beta} \varepsilon(s) + \varepsilon'(s) \right) \dot{x}(s) ds \\ &\quad - \alpha(t)x(t) + e^{-\frac{t}{\beta}} \alpha_0 x_0 + \int_0^t e^{\frac{s-t}{\beta}} \left(\alpha'(s) + \frac{\alpha(s)}{\beta} \right) x(s) ds. \end{aligned}$$

One then proves (3.8) by inserting (C.3) in (3.3) and by using the uniform boundedness of x and Assumption 3. The proof of (3.9) follows the same arguments based on the fact that $\forall t \geq 0$,

$$\nabla f(x_{LM}(t)) = e^{-\frac{t}{\beta}} \nabla f(x_0) - e^{-\frac{t}{\beta}} \int_0^t \frac{1}{\beta} e^{\frac{s}{\beta}} \alpha(s) \dot{x}_{LM}(s) ds.$$

Appendix D. Integrability of φ and Additional Asymptotic Computations. Below we prove Lemma 4.7.

Proof. We suppose that Assumptions 4 and 5 hold. As stated in Remark 4.3, since φ is continuous, we only need to check its integrability when t tends to $+\infty$. Let $t > 0$, we first establish some useful identities, we omit the dependence on t for the sake of readability.

$$p' = \frac{\alpha' \varepsilon - (\alpha + \beta \lambda) \varepsilon'}{\varepsilon^2}, \quad \text{and} \quad p'' = \frac{\alpha'' \varepsilon^2 - 2\alpha' \varepsilon' \varepsilon - (\alpha + \beta \lambda) \varepsilon'' \varepsilon + 2(\alpha + \beta \lambda) (\varepsilon')^2}{\varepsilon^3}.$$

Then,

$$\begin{aligned} (D.1) \quad r &= \frac{p^2}{4} \left(1 + \frac{2p'}{p^2} - \frac{4\lambda}{\varepsilon p^2} \right) = \frac{(\alpha + \beta \lambda)^2}{4\varepsilon^2} \left(1 + \frac{2p' \varepsilon^2}{(\alpha + \beta \lambda)^2} - \frac{4\lambda \varepsilon}{(\alpha + \beta \lambda)^2} \right) \\ &= \frac{(\alpha + \beta \lambda)^2}{4\varepsilon^2} \left(1 + \frac{2\alpha' \varepsilon}{(\alpha + \beta \lambda)^2} - \frac{2\varepsilon'}{(\alpha + \beta \lambda)} - \frac{4\lambda \varepsilon}{(\alpha + \beta \lambda)^2} \right). \end{aligned}$$

An important consequence of Assumption 5 is that $|\varepsilon'(t)| = o(\varepsilon(t))$, $|\varepsilon''(t)| = o(\varepsilon'(t))$ (and the same holds for α w.r.t. to its derivatives). Therefore, we deduce from (D.1) that

$$r(t) \sim_{+\infty} \frac{(\alpha(t) + \beta \lambda)^2}{4\varepsilon(t)^2},$$

and we note that $1/r$ decays at the same speed as ε^2 , which will be useful later. In order to study φ , we now differentiate r ,

$$\begin{aligned} (D.2) \quad r' &= \frac{p' p}{2} \left(1 + \frac{2p'}{p^2} - \frac{4\lambda}{\varepsilon p^2} \right) + \frac{1}{4} \left(2p'' - \frac{4(p')^2}{p} + \frac{8\lambda p'}{\varepsilon p} + \frac{4\lambda \varepsilon'}{\varepsilon^2} \right) \\ &= \frac{2p'}{p} r + \frac{1}{4} \left(2p'' - \frac{4(p')^2}{p} + \frac{8\lambda p'}{\varepsilon p} + \frac{4\lambda \varepsilon'}{\varepsilon^2} \right), \quad \text{and,} \end{aligned}$$

$$(D.3) \quad r'' = 2\frac{p''p - (p')^2}{p^2}r + \frac{2p'}{p}r' \\ + \frac{1}{4} \left(2p''' + 4\frac{(p')^3 - 2p''p'}{p^2} + 8\lambda\frac{p''p\varepsilon - (p')^2\varepsilon - p'p\varepsilon'}{\varepsilon^2p^2} + \frac{4\lambda\varepsilon''}{\varepsilon^2} - \frac{8\lambda(\varepsilon')^2}{\varepsilon^3} \right).$$

Then, to justify that φ is integrable, we prove that $\frac{r''}{r^{3/2}}$ and $\frac{(r')^2}{r^{5/2}}$ are integrable. Since we know that $1/r$ decays at the same speed as ε^2 , we can equivalently show that $\varepsilon^3 r''$ and $\varepsilon^5 (r')^2$ are integrable. To this aim we fully expand all the terms in (D.2) and (D.3). We first deal with (D.2), it holds that:

$$r'^2 \varepsilon^5 = \left[\frac{(\alpha + \beta\lambda)^2 \left(-\frac{4\lambda\varepsilon}{(\alpha + \beta\lambda)^2} + 1 + \frac{(-2(\alpha + \beta\lambda)\varepsilon' + 2\alpha'\varepsilon)}{(\alpha + \beta\lambda)^2} \right) \varepsilon'}{2\sqrt{\varepsilon}} \right. \\ + \frac{(\alpha + \beta\lambda)^2 \sqrt{\varepsilon}}{4} \left(-\frac{4\lambda\varepsilon'}{(\alpha + \beta\lambda)^2} + \frac{8\lambda\alpha'\varepsilon}{(\alpha + \beta\lambda)^3} + \frac{2 \left(-\frac{2(\alpha + \beta\lambda)\varepsilon'}{\varepsilon} + 2\alpha'\varepsilon \right) \varepsilon'}{(\alpha + \beta\lambda)^2} \right. \\ \left. \left. + \frac{\left(\frac{4(\alpha + \beta\lambda)\varepsilon'^2}{\varepsilon} - 2(\alpha + \beta\lambda)\varepsilon'' - 4\alpha'\varepsilon' + 2\alpha''\varepsilon \right)}{(\alpha + \beta\lambda)^2} - \frac{2(-2(\alpha + \beta\lambda)\varepsilon' + 2\alpha'\varepsilon)\alpha'}{(\alpha + \beta\lambda)^3} \right) \right. \\ \left. + \frac{(\alpha + \beta\lambda)}{2} \left(-\frac{4\lambda\varepsilon}{(\alpha + \beta\lambda)^2} + 1 + \frac{(-2(\alpha + \beta\lambda)\varepsilon' + 2\alpha'\varepsilon)}{(\alpha + \beta\lambda)^2} \right) \alpha' \sqrt{\varepsilon} \right]^2.$$

The computations for $\varepsilon^2 r''$ are omitted since they are longer but very similar.

We analyze the integrability of each of the terms above (and those of $\varepsilon^2 r''$). By Assumption 5, ε' , ε'' and ε''' are integrable, as well as α' , α'' and α''' , this justifies the integrability of most of the terms. We also need $\frac{(\varepsilon')^2}{\varepsilon}$ and $\frac{(\varepsilon')^3}{\varepsilon}$ to be integrable, which holds by Assumption 5. Overall, φ is integrable on \mathbb{R}_+ . \square

We now state and prove the a result used at the end of the proof of Theorem 4.8.

LEMMA D.1. *Under Assumptions 4 and 5, for all $s \geq 0$,*

$$\frac{1}{16} \left(\frac{2p'(s)}{p(s)^{3/2}} - \frac{4\lambda\sqrt{\varepsilon(s)}}{(\alpha(s) + \beta\lambda)^{3/2}} \right)^2 = \frac{\lambda^2\varepsilon(s)}{(\alpha(s) + \beta\lambda)^3} + o(\varepsilon(s)).$$

Proof. We omit the time dependence on $s \geq 0$ for the sake of readability. Using Assumption 4 we can define and expand the following quantity,

$$\frac{1}{16} \left(\frac{2p'}{p^{3/2}} - \frac{4\lambda\sqrt{\varepsilon}}{(\alpha + \beta\lambda)^{3/2}} \right)^2 = \frac{\lambda^2\varepsilon}{(\alpha + \beta\lambda)^3} - \frac{p'\lambda\sqrt{\varepsilon}}{p^{3/2}(\alpha + \beta\lambda)^{3/2}} + \frac{(p')^2}{4p^3} \\ = \frac{\lambda^2\varepsilon}{(\alpha + \beta\lambda)^3} - \frac{\lambda(\alpha'\varepsilon - (\alpha + \beta\lambda)\varepsilon')}{(\alpha + \beta\lambda)^3} + \frac{(\alpha')^2\varepsilon + \frac{(\varepsilon')^2}{\varepsilon}(\alpha + \beta\lambda)^2 - 2\alpha'\varepsilon'(\alpha + \beta\lambda)}{4(\alpha + \beta\lambda)^3}.$$

Assumption 5, implies in particular that $|\varepsilon'(t)| = o(\varepsilon(t))$ and that $\alpha'(t) \rightarrow 0$, which we use in the equality above to obtain the desired conclusion. \square

Appendix E. Additional Experimental Details. We used Euler discretization schemes with fixed step-size $\gamma > 0$, for approximating the solutions of the three ODEs considered in Section 5. For a trajectory x , at times $t_k = \gamma k$, $k \in \mathbb{N}$, we denote $x(t_k) \stackrel{\text{def}}{=} x^{(k)}$. We approximated (CN) by explicit discretization, $\forall k \in \mathbb{N}$,

$$(E.1) \quad x_N^{(k+1)} = x_N^{(k)} - \gamma \left[\beta \nabla^2 f(x_N^{(k)}) \right]^{-1} \nabla f(x_N^{(k)}).$$

Then, defining $\varepsilon_k = \varepsilon(t_k)$ and $\alpha_k = \alpha(t_k)$, (LM) and (VM-DIN-AVD) are obtained via Euler semi-implicit discretization. The solution of (LM) is approximated by,

$$(E.2) \quad x_{LM}^{(k+1)} = x_{LM}^{(k)} - \gamma \left[\alpha_k I_n + \beta \nabla^2 f(x_{LM}^{(k)}) \right]^{-1} \nabla f(x_{LM}^{(k)}),$$

where I_n is the identity matrix on \mathbb{R}^n . The solution of (VM-DIN-AVD) is similarly,

$$(E.3) \quad x^{(k+1)} = x^{(k)} + \left[(\varepsilon_k + \gamma \alpha_k) I_n + \gamma \beta \nabla^2 f(x^{(k)}) \right]^{-1} \left(\varepsilon_k (x^{(k)} - x^{(k-1)}) - \gamma^2 \nabla f(x^{(k)}) \right).$$

As sanity check, one can see that for $\varepsilon_k = 0$, (E.3) is equivalent to (E.2), which is itself equivalent to (E.1) when $\alpha_k = 0$.

References.

- [1] C. D. ALECSA, S. C. LÁSZLÓ, AND T. PINȚA, *An extension of the second order dynamical system that models Nesterov's convex gradient method*, Applied Mathematics & Optimization, 84 (2021), pp. 1687–1716.
- [2] F. ALVAREZ, H. ATTOUCH, J. BOLTE, AND P. REDONT, *A second-order gradient-like dissipative dynamical system with Hessian-driven damping: Application to optimization and mechanics*, Journal de Mathématiques Pures et Appliquées, 81 (2002), pp. 747–779.
- [3] F. ALVAREZ AND J. PÉREZ, *A dynamical system associated with Newton's method for parametric approximations of convex minimization problems*, Applied Mathematics and Optimization, 38 (1998), pp. 193–217.
- [4] H. ATTOUCH, M. M. ALVES, AND B. F. SVAITER, *A dynamic approach to a proximal-Newton method for monotone inclusions in Hilbert spaces, with complexity $O(1/n^2)$* , Journal of Convex Analysis, 23 (2016), pp. 139–180.
- [5] H. ATTOUCH, A. BALHAG, Z. CHBANI, AND H. RIAHI, *Fast convex optimization via inertial dynamics combining viscous and Hessian-driven damping with time rescaling*, Evolution Equations and Control Theory, 11 (2022), pp. 487–514.
- [6] H. ATTOUCH, R. I. BOȚ, AND E. R. CSETNEK, *Fast optimization via inertial dynamics with closed-loop damping*, Journal of the European Mathematical Society (published online), (2022).
- [7] H. ATTOUCH AND A. CABOT, *Asymptotic stabilization of inertial gradient dynamics with time-dependent viscosity*, Journal of Differential Equations, 263 (2017), pp. 5412–5458.
- [8] H. ATTOUCH, Z. CHBANI, J. FADILI, AND H. RIAHI, *First-order optimization algorithms via inertial systems with Hessian driven damping*, Math. Program., 194 (2020), pp. 1–43.
- [9] H. ATTOUCH AND J. FADILI, *From the ravine method to the Nesterov method and vice versa: A dynamical system perspective*, SIAM Journal on Optimization, 32 (2022), pp. 2074–2101.
- [10] H. ATTOUCH AND S. C. LÁSZLÓ, *Newton-like inertial dynamics and proximal algorithms governed by maximally monotone operators*, SIAM Journal on Optimization, 30 (2020), pp. 3252–3283.

- [11] H. ATTOUCH AND S. C. LÁSZLÓ, *Continuous Newton-like inertial dynamics for monotone inclusions*, Set-Valued and Variational Analysis, 29 (2021), pp. 555–581.
- [12] H. ATTOUCH, J. PEYPOUQUET, AND P. REDONT, *Fast convex optimization via inertial dynamics with Hessian driven damping*, Journal of Differential Equations, 261 (2016), pp. 5734–5783.
- [13] H. ATTOUCH, P. REDONT, AND B. F. SVAITER, *Global convergence of a closed-loop regularized Newton method for solving monotone inclusions in hilbert spaces*, Journal of Optimization Theory and Applications, 157 (2013), pp. 624–650.
- [14] H. ATTOUCH AND B. F. SVAITER, *A continuous dynamical Newton-like approach to solving monotone inclusions*, SIAM Journal on Control and Optimization, 49 (2011), pp. 574–598.
- [15] H. BAUSCHKE AND P. COMBETTES, *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*, Springer, 2011.
- [16] M. BENAÏM, J. HOFBAUER, AND S. SORIN, *Stochastic approximations and differential inclusions*, SIAM Journal on Control and Optimization, 44 (2005), pp. 328–348.
- [17] C. M. BENDER AND S. A. ORSZAG, *Asymptotic methods and perturbation theory*, Springer, 1999.
- [18] O. BLUMENTHAL, *Über asymptotische Integration linearer Differentialgleichungen mit Anwendung auf eine asymptotische Theorie der Kugelfunktionen*, Archiv der Mathematik und Physik, 19 (1912), pp. 136–174.
- [19] R. I. BOŦ, E. R. CSETNEK, AND S. C. LÁSZLÓ, *Tikhonov regularization of a second order dynamical system with Hessian driven damping*, Math. Program., 189 (2021), pp. 151–186.
- [20] H. BREZIS, *Opérateurs maximaux monotones et semi-groupes de contractions dans les espaces de Hilbert*, Elsevier, 1973.
- [21] L. BRILLOUIN, *Remarques sur la mécanique ondulatoire*, Journal de Physique et Le Radium, 7 (1926), pp. 353–368.
- [22] C. G. BROYDEN, *The convergence of a class of double-rank minimization algorithms*, IMA Journal of Applied Mathematics, 6 (1970), pp. 76–90.
- [23] C. CASTERA, *Inertial Newton algorithms avoiding strict saddle points*, arXiv:2111.04596, (2021).
- [24] C. CASTERA, J. BOLTE, C. FÉVOTTE, AND E. PAUWELS, *An inertial Newton algorithm for deep learning*, Journal of Machine Learning Research, 22 (2021), pp. 1–31.
- [25] L. CHEN AND H. LUO, *First order optimization methods based on Hessian-driven Nesterov accelerated gradient flow*, arXiv:1912.09276, (2019).
- [26] D. DRUSVYATSKIY, A. D. IOFFE, AND A. S. LEWIS, *Nonsmooth optimization using Taylor-like models: error bounds, convergence, and termination criteria*, Mathematical Programming, 185 (2021), pp. 357–383.
- [27] E. EMMRICH, *Discrete versions of Gronwall’s lemma and their application to the numerical analysis of parabolic problems*, TU Fachbereich, 1999.
- [28] R. FLETCHER, *A new approach to variable metric algorithms*, The computer journal, 13 (1970), pp. 317–322.
- [29] M. K. GAVURIN, *Nonlinear functional equations and continuous analogues of iteration methods*, Izvestiya Vysshikh Uchebnykh Zavedenii. Matematika, 1 (1958), pp. 18–31.
- [30] D. GOLDFARB, *A family of variable-metric methods derived by variational means*, Mathematics of computation, 24 (1970), pp. 23–26.

- [31] G. GREEN, *On the motion of waves in a variable canal of small depth and width*, Transactions of the Cambridge Philosophical Society, 6 (1838), pp. 457–462.
- [32] J. D. HUNTER, *Matplotlib: A 2D graphics environment*, Computing in science & engineering, 9 (2007), pp. 90–95.
- [33] H. A. KRAMERS, *Wellenmechanik und halbzahlige Quantisierung*, Zeitschrift für Physik, 39 (1926), pp. 828–840.
- [34] T. LIN AND M. I. JORDAN, *A control-theoretic perspective on optimal high-order optimization*, Math. Program., 195 (2022), pp. 929–975.
- [35] J. LIOUVILLE, *Mémoire sur le développement des fonctions ou parties de fonctions en séries dont les divers termes sont assujétis à satisfaire à une même équation différentielle du second ordre, contenant un paramètre variable.*, Journal de Mathématiques Pures et Appliquées, 2 (1836), pp. 253–265.
- [36] D. C. LIU AND J. NOCEDAL, *On the limited memory BFGS method for large scale optimization*, Math. Program., 45 (1989), pp. 503–528.
- [37] L. LJUNG, *Analysis of recursive stochastic algorithms*, IEEE transactions on automatic control, 22 (1977), pp. 551–575.
- [38] Y. NESTEROV, *A method for unconstrained convex minimization problem with the rate of convergence $O(1/k^2)$* , in Doklady an USSR, vol. 269, 1983, pp. 543–547.
- [39] Y. NESTEROV, *Introductory lectures on convex optimization: A basic course*, vol. 87, Springer Science & Business Media, 2003.
- [40] P. OCHS, J. FADILI, AND T. BROX, *Non-smooth non-convex Bregman minimization: Unification and new algorithms*, Journal of Optimization Theory and Applications, 181 (2019), pp. 244–278.
- [41] F. OLVER, *Error bounds for the Liouville–Green (or WKB) approximation*, in Mathematical Proceedings of the Cambridge Philosophical Society, vol. 57, Cambridge University Press, 1961, pp. 790–810.
- [42] F. OLVER, *Asymptotics and special functions*, Academic Press, 1997.
- [43] B. T. POLYAK, *Some methods of speeding up the convergence of iteration methods*, USSR Computational Mathematics and Mathematical Physics, 4 (1964), pp. 1–17.
- [44] G. ROSSUM, *Python reference manual*, CWI (Centre for Mathematics and Computer Science), 1995.
- [45] D. F. SHANNO, *Conditioning of quasi-Newton methods for function minimization*, Mathematics of computation, 24 (1970), pp. 647–656.
- [46] B. SHI, S. S. DU, M. I. JORDAN, AND W. J. SU, *Understanding the acceleration phenomenon via high-resolution differential equations*, Math. Program., 195 (2022), pp. 79–148.
- [47] W. SU, S. BOYD, AND E. CANDÈS, *A differential equation for modeling Nesterov’s accelerated gradient method: Theory and insights*, in Advances in Neural Information Processing Systems (NeurIPS), vol. 27, 2014, pp. 2510–2518.
- [48] J. G. TAYLOR, *Improved error bounds for the Liouville–Green (or WKB) approximation*, Journal of Mathematical Analysis and Applications, 85 (1982), pp. 79–89.
- [49] S. V. D. WALT, C. COLBERT, AND G. VAROQUAUX, *The NumPy array: a structure for efficient numerical computation*, Computing in Science & Engineering, 13 (2011), pp. 22–30.
- [50] G. WENTZEL, *Eine Verallgemeinerung der Quantenbedingungen für die Zwecke der Wellenmechanik*, Zeitschrift für Physik, 38 (1926), pp. 518–529.