

# INNER PRODUCT PRECONDITIONED TRUST-REGION METHODS FOR FREQUENCY-DOMAIN FULL WAVEFORM INVERSION\*

XAVIER ADRIAENS<sup>†</sup>, LUDOVIC MÉTIVIER<sup>‡</sup>, AND CHRISTOPHE GEUZAINÉ<sup>§</sup>

**Abstract.** Full waveform inversion is a seismic imaging method which requires to solve a large-scale minimization problem, typically through local optimization techniques. Most local optimization methods can basically be built up from two choices: the update direction and the strategy to control its length. In the context of full waveform inversion, this strategy is very often a line search. We here propose to use instead a trust-region method, in combination with non-standard inner products which act as preconditioners. More specifically, a line search and several trust-region variants of the steepest descent, the limited memory BFGS algorithm and the inexact Newton method are presented and compared. A strong emphasis is given to the inner product choice. For example, its link with preconditioning the update direction and its implication in the trust-region constraint are highlighted. A first numerical test is performed on a 2D synthetic model then a second configuration, containing two close reflectors, is studied. The latter configuration is known to be challenging because of multiple reflections. Based on these two case studies, the importance of an appropriate inner product choice is highlighted and the best trust-region method is selected and compared to the line search method. In particular we were able to demonstrate that using an appropriate inner product greatly improves the convergence of all the presented methods and that inexact Newton methods should be combined with trust-region methods to increase their convergence speed.

**Key words.** numerical optimization, large-scale inverse problems, trust-regions methods, operator preconditioning, seismic imaging, full waveform inversion.

**AMS subject classifications.** 35R30, 65K10, 86-08, 49M15, 90C06

**1. Introduction.** Full waveform inversion is a high-resolution seismic imaging technique formulated as a data fitting problem, whose aim is to recover some model parameters by minimizing the discrepancy between recorded data and data simulated by solving wave propagation problems [29, 35]. By nature these data are oscillatory and consequently the misfit quantifying the discrepancy features local minima [4, 21]. Global optimization techniques should ideally be used but the typically very high dimensions of the search space prohibits their use and only local optimization tools can practically be employed, with care [8]. A straightforward direction to iteratively update the model properties is of course the gradient, *i.e.* the direction of steepest decrease. However it is well-known that the inverse Hessian plays a crucial role in the reconstruction in addition to offering the possibility to account for coupling effects between parameter classes for multi-parameter inversion [3, 24, 26, 29, 38]. A theoretically simple way to incorporate these second-order derivatives is to minimize the misfit using Newton methods. In practice however the pure Newton method is too computationally intensive to implement, because it requires inverting the Hessian operator. In addition, the misfit is not necessarily quadratic, thus the exact Newton direction is not necessarily appropriate. Consequently, it is natural to turn to inexact Newton methods, where the search direction is constructed iteratively to approximate the pure Newton direction, or to quasi-Newton methods. State-of-the-art methods rely on the quasi-Newton *l*-BFGS algorithm, which implicitly builds an approximation of the inverse Hessian operator from *l* previously saved gradients and model parameters [23]. However it has been illustrated that on some specific cases involving multiple re-

---

\*Submitted to the editors March 8, 2022.

**Funding:** This research was funded by the Fonds de la Recherche Scientifique de Belgique (F.R.S.-FNRS) and by the ARC “WAVES” grant 15/19-03 from the Wallonia-Brussels Federation of Belgium. Computational resources were provided by the Consortium des Équipements de Calcul Intensif (CÉCI), funded by the F.R.S.-FNRS and by the Walloon Region.

<sup>†</sup>F.R.S.-FNRS, Department of Electrical Engineering and Computer Science, Université de Liège, Belgium ([xavier.adriaens@uliege.be](mailto:xavier.adriaens@uliege.be)).

<sup>‡</sup>Univ. Grenoble Alpes, CNRS, LJK, F-38000 Grenoble, France ([ludovic.metivier@univ-grenoble-alpes.fr](mailto:ludovic.metivier@univ-grenoble-alpes.fr)).

<sup>§</sup>Department of Electrical Engineering and Computer Science, Université de Liège, Belgium ([cgeuzaine@uliege.be](mailto:cgeuzaine@uliege.be)).

45 flections, such quasi-Newton methods fail to converge where inexact Newton methods  
 46 do succeed [20]. The latter compute the descent direction through a few iterations of  
 47 a linear system involving the Hessian operator (the Newton system). One advantage  
 48 over  $l$ -BFGS is the locality of the quadratic approximation: such methods do not  
 49 rely on the convergence history of the algorithm, which might yield inaccurate inverse  
 50 Hessian approximation for non quadratic misfit functions. The bottleneck of these  
 51 methods lies in the compromise to find between a direction built in few iterations, but  
 52 which hardly takes the Hessian into account and a nearly exact direction which is very  
 53 expensive to compute. A complementary strategy to reduce this number of inner it-  
 54 eration is to apply a preconditioner to both sides of the Newton system [6, 16, 25, 36].  
 55 To implement any of the three above mentioned schemes, one can rely either on line  
 56 search algorithms, or on trust-region methods. In the former case, once a direction  
 57 is chosen, the outer iteration is completed by finding the optimal length of the step  
 58 that should be performed along that direction. Among the non linear optimization  
 59 community, it is sometimes argued however that line search is not well suited with  
 60 Newton directions, especially when the Hessian is nearly singular. Indeed when the  
 61 Hessian is nearly singular, the Newton direction becomes excessively long such that  
 62 the local quadratic approximation implicitly made when computing it ceases to hold.  
 63 Much computational effort must then be made by the line search procedure to reduce  
 64 the step size [23]. Stopping the iterative solution of the Newton system earlier ap-  
 65 pears as a solution to this problem. For example, its convergence requirements could  
 66 be relaxed such that they reflect the accuracy of the local quadratic approximation  
 67 [9, 19]. Alternatively, a trust-region method could be used instead [18, 37, 39, 40].  
 68 The latter limits the length of the update direction depending on the accuracy of the  
 69 local quadratic approximation. The length of a direction is given by its norm, itself  
 70 induced by the inner product chosen for the model parameters space. The choice of  
 71 this inner product is thus pivotal in the implementation of a trust-region method.  
 72 Moreover changing the inner product modifies both the gradient and the Hessian and  
 73 is equivalent to applying a preconditionner [7, 12, 22, 17, 41]. Consequently it also  
 74 has a major impact on line search based local optimization methods.

75 In this paper, we tackle the three following important questions:

- 76 • Which descent direction to compute: the gradient, the  $l$ -BFGS direction or  
 77 an inexact Newton direction?
- 78 • Which globalization method to select: a line search method or a trust-region  
 79 method?
- 80 • Which preconditioning strategy to apply? How to enforce it?

81 Answering these three questions and determining the good combinations (good prac-  
 82 tices) between them is crucial for effective full waveform inversion. From our study, it  
 83 appears that preconditioning is essential and that enforcing preconditioning through  
 84 the inner product is elegant and, more interestingly, implies no modification to the  
 85 practical implementation of the optimization algorithms. The  $l$ -BFGS method is  
 86 found to be the most efficient method for the considered single-parameter inversions.  
 87 It is also found to be insensitive to the globalization choice. Inexact Newton methods  
 88 should not be discarded though, as considering the exact Hessian might lead to better  
 89 model parameter decoupling in the case multi-parameter inversions. When using in-  
 90 exact Newton methods, our case studies show that using a trust region globalization  
 91 consistently improves convergence.

92 The paper is organized as follows. In the first part, full waveform inversion is  
 93 stated very generally. The optimization problem and its solution procedures using  
 94 either a line search or a trust-region are introduced. The Newton system, which is  
 95 pivotal in local minimization theory, is also derived. A particular emphasis is given to  
 96 the inner product choice. More specifically, its link with preconditioning the Newton  
 97 system is established. Local minimization methods commonly used in the context  
 98 of full waveform inversion are then recalled. In the second part, the application

99 to acoustic imaging is detailed. The (adjoint) procedure to compute gradients and  
 100 Hessian vector products is given and its computational cost is explained. The overall  
 101 computational cost of each optimization method is then deduced. Finally, convergence  
 102 results on the acoustic Marmousi case study are analyzed to determine the best inner  
 103 product and the best parameters for a trust-region method. This best candidate is  
 104 then compared to line search methods on both the Marmousi model and on a case  
 105 study involving strong reflectors.

106 **2. Local optimization methods.** Full wave inversion is an imaging method  
 107 based on the minimization of a misfit functional  $J$ , which exclusively depends on  
 108 some model parameters  $m$ . The recovered model parameters  $m^*$  are defined as the  
 109 minimizer of this misfit, *i.e.*  $m^* = \arg \min J(m)$ . Local optimization techniques are  
 110 based on a local quadratic expansion of the misfit  $J$  around the current model estimate

$$111 \quad (2.1) \quad J(m + \delta m) \approx J(m) + \{D_m J\}(\delta m) + \frac{1}{2} \{D_{mm}^2 J\}(\delta m, \delta m).$$

112 This expansion can also be written in terms of the gradient  $j'$  and the Hessian operator  
 113  $H$  once an inner product  $\langle \cdot, \cdot \rangle_M$  is chosen for the model space  $M$

$$114 \quad (2.2) \quad J(m + \delta m) \approx J(m) + \langle j', \delta m \rangle_M + \frac{1}{2} \langle H \delta m, \delta m \rangle_M.$$

115 The pure Newton direction  $p_N$  is then defined as the minimizer of this local quadratic  
 116 expansion, which is also the solution of a linear system

$$117 \quad (2.3) \quad p_N = \arg \min_{p \in M} J(m) + \langle j', p \rangle_M + \frac{1}{2} \langle H p, p \rangle_M \quad \text{or} \quad H p_N = -j'.$$

118 The large-scale nature of this linear system requires either the use of approximate  
 119 Hessian operators that are straightforward to invert, or the use of Hessian-free iterative  
 120 methods. Both approaches are usually referred to as quasi-Newton methods and  
 121 inexact Newton methods. In the latter case, the conjugate gradient method is the  
 122 ideal candidate for the iterative solver because the Hessian operator is symmetric. The  
 123 conjugate gradient method is however designed for positive definite operators while  
 124 the full Hessian can be indefinite, especially far from the global minimum [29, 20].  
 125 As a consequence, either an additional safeguard is added to exit prematurely when  
 126 directions of negative curvature are encountered or the exact Hessian is modified such  
 127 that it becomes positive definite, *e.g.* using the Gauss-Newton approximation [25].

128 **2.1. Globalization methods.** As mentioned in the introduction, the misfit is  
 129 not necessarily quadratic and thus the pure Newton direction or its approximations  
 130 are not always the best directions. For that reason the length of the search direction is  
 131 often tweaked using a line search or a trust-region method, which ensures convergence  
 132 towards the nearest local minimum [9, 11, 10, 23].

133 **2.1.1. Line search.** When using a line search procedure, a direction  $p$  must first  
 134 be identified. An appropriate length  $\gamma$  is then given to this direction  $p$ , ideally the  
 135 global minimum along the line  $m + \gamma p$ . In practice however less stringent satisfactory  
 136 conditions are used instead to spare expensive wave problem resolutions. Maybe the  
 137 best example are strong Wolfe conditions

$$138 \quad (2.4) \quad J(m + \gamma p) \leq J(m) + c_1 \gamma \{D_m J(m)\}(p)$$

$$139 \quad (2.5) \quad |\{D_m J(m + \gamma p)\}(p)| \leq c_2 |\{D_m J(m)\}(p)|$$

141 for some constant  $c_1$  and  $c_2$  such that  $0 < c_1 < c_2 < 1$ . The first condition is called the  
 142 sufficient decrease condition. It ensures that updating the model in the direction  $\gamma p$

143 produces a decrease smaller than a fraction  $c_1$  of what is expected from a local linear  
 144 approximation of the misfit. The second condition, called the curvature condition,  
 145 ensures that the updated model  $m + \gamma p$  is sufficiently close to a local minimum along  
 146 the line, where the directional derivative  $\{D_m J(m + \gamma p)\}(p)$  would be zero. When this  
 147 derivative is very smaller (resp. larger) than zero, then a larger (resp. smaller) step  
 148 could produce a significantly bigger decrease. We choose here a line search algorithm  
 149 that satisfies strong Wolfe conditions and accepts steps easily (Algorithm 3.2 from  
 150 [23] with  $c_1 = 10^{-4}$ ,  $c_2 = 0.9$ ). The outer loop is finally obtained by repeating these  
 151 two steps iteratively until convergence.

152 **2.1.2. Trust region.** At the opposite when using a trust-region method, first  
 153 a maximum length  $\Delta$  is chosen. Then the best approximate solution, meaning the  
 154 direction that minimizes a local prediction of the misfit but smaller than this length,  
 155 is used

$$156 \quad (2.6) \quad p = \arg \min_{p \in M, \|p\|_M \leq \Delta} \left[ J^{\text{pred}}(m; p) := J(m) + \langle j'(m), p \rangle_M + 0.5 \langle \tilde{H}(m)p, p \rangle_M \right].$$

157 This local misfit prediction  $J^{\text{pred}}$  is typically constructed based on the local quadratic  
 158 approximation (2.2) through a particular choice of some approximate Hessian operator  
 159  $\tilde{H}$ . Of course the approximate Newton direction  $\tilde{H}p = -j'$  is the solution of this  
 160 problem if it lies inside the trust region. There are several possibilities to choose  
 161 this length  $\Delta$  and our particular choice is detailed later. More importantly, as we  
 162 pointed out in the introduction, the length constraint is formulated in terms of the  
 163 norm induced by the inner product  $\|p\|_M^2 = \langle p, p \rangle_M \leq \Delta^2$ . Modifying this inner  
 164 product therefore changes the shape of the trust region and it is then desirable to  
 165 choose it carefully. The size of the trust region is actually controlled by the outer  
 166 iterations. The decision of modifying the trust region is based on the accuracy of the  
 167 local prediction of the misfit. When the prediction is accurate but the updates are  
 168 limited by the length constraint, then the trust region radius is increased. At the  
 169 opposite, when the updates are out of the range of validity of the prediction, then  
 170 the trust region radius is decreased. The decrease (resp. increase) rate of the radius  
 171 is controlled by some parameter  $c_0 < 1$  (resp.  $c_1 > 1$ ). The quality of the prediction  
 172 is quantified by the ratio between the actual decrease  $\delta J_a := J(m_n) - J(m_{n+1})$  and  
 173 the decrease predicted by the local prediction of the misfit. There are two ways to  
 174 compute this predicted decrease [10]. On the one hand the expansion can be written  
 175 in terms of the gradient and the Hessian operator at the previous model estimate

$$176 \quad (2.7) \quad J(m_{n+1}) = J(m_n + p_n)$$

$$177 \quad (2.8) \quad \approx J(m_n) + \langle j'(m_n), p_n \rangle_M + 0.5 \langle \tilde{H}(m_n)p_n, p_n \rangle_M = J^{\text{pred}}(m_n; p_n)$$

179 which defines the prospective predicted decrease

$$180 \quad (2.9) \quad \delta J_{p,p} := J(m_n) - J^{\text{pred}}(m_n; p_n)$$

$$181 \quad (2.10) \quad = -\langle j'(m_n), p_n \rangle_M - 0.5 \langle \tilde{H}(m_n)p_n, p_n \rangle_M.$$

183 On the other hand, it can also be written in terms of the gradient and the Hessian  
 184 operator at the next model estimate

$$185 \quad J(m_n) = J(m_{n+1} - p_n)$$

$$186 \quad \approx J(m_{n+1}) - \langle j'(m_{n+1}), p_n \rangle_M + 0.5 \langle \tilde{H}(m_{n+1})p_n, p_n \rangle_M = J^{\text{pred}}(m_{n+1}; -p_n)$$

187

188 which defines the retrospective predicted decrease

$$189 \quad (2.11) \quad \delta J_{p,r} := J^{\text{pred}}(m_{n+1}; -p_n) - J(m_{n+1})$$

$$190 \quad (2.12) \quad = -\langle j'(m_{n+1}), p_n \rangle_M + 0.5 \left\langle \tilde{H}(m_{n+1}) p_n, p_n \right\rangle_M.$$

192 These ratios between the actual decrease and one of both the predicted decreases  $\rho_p :=$   
 193  $\delta J_a / \delta J_{p,p}$  and  $\rho_r := \delta J_a / \delta J_{p,r}$  are actually both equal to one when the approximate  
 194 Hessian in the update direction and the second order expansion (2.2) are exact. When  
 195 the misfit is not quadratic or the Hessian approximation is not accurate, then these  
 196 ratios can go away from one. Using anything else than the full Newton method can  
 197 degrade these ratios, even if the misfit is quadratic. In particular for a pure quadratic  
 198 misfit, neglecting the negative definite part of the Hessian makes the prospective ratio  
 199 bigger than one ( $\delta J_{p,p}$  is underestimated) and the retrospective ratio smaller than one  
 200 ( $\delta J_{p,r}$  is overestimated).

201 Standard trust-region methods directly control the radius  $\Delta$ . However it is an  
 202 absolute quantity, in the sense that it is compared to  $\|p\|_M$ , which depends on the  
 203 inner product. Thus, it seems more natural to control this radius relatively to the  
 204 gradient norm ( $\Delta = \mu \|j'\|_M$ ), which provides a length reference for the (approximate)  
 205 Newton system. In this way, even when the (approximate) Newton system changes  
 206 scale from one iteration to another, the trust region remains relevant. This particular  
 207 variant (Algorithm 2.1) has been first introduced in [11].

---

**Algorithm 2.1** Fan trust-region algorithm
 

---

**Require:** retrospective or prospective,  $0 \leq \rho_0 < \rho_1 < 1$  and  $0 < c_0 < 1 < c_1$

$\mu_0 = 1$

**loop**

$\Delta_n = \mu_n \|j'(m_n)\|_M$

$p_n = \begin{cases} -\mu_n j'_n \\ (2.28) \text{ with } \Delta = \Delta_n \\ \text{Algorithm 2.5 with } \Delta = \Delta_n \end{cases}$

$\delta J_a = J(m_n) - J(m_n + p_n)$  and  $\delta J_{p,p} = J(m_n) - J^{\text{pred}}(m_n; p_n)$

$\rho_p = \delta J_a / \delta J_{p,p}$

**if**  $\rho_p \geq \rho_0$  **then**  $m_{n+1} = m_n + p_n$  **else**  $m_{n+1} = m_n$

**if** prospective **or**  $\rho_p < \rho_0$  **then**

$\rho = \rho_p$

**else if** retrospective **then**

$\delta J_{p,r} = J^{\text{pred}}(m_{n+1}; -p_n) - J(m_{n+1})$

$\rho = \rho_r = \delta J_a / \delta J_{p,r}$

**end if**

**if**  $\rho < \rho_1$

**then**  $\mu_{n+1} = c_0 \mu_n$

**else if**  $\rho \geq \rho_1$  **and**  $\|p_n\|_M > 0.5 \Delta_n$

**then**  $\mu_{n+1} = c_1 \mu_n$

**else**

**then**  $\mu_{n+1} = \mu_n$

**end loop**

---

208 According to this algorithm, a direction  $p_n$  is rejected when the prospective misfit  
 209 prediction  $J_n^{\text{pred}}$  used to compute it is not accurate, in the sense that the prospective  
 210 ratio is smaller than some threshold  $\rho_0$ . If not rejected, then the trust region size  
 211 is updated according to either the prospective or the retrospective ratio, based on a  
 212 comparison with a second threshold  $\rho_1$ . Because the updated radius  $\Delta_{n+1}$  constrains  
 213 the direction search around the next model estimate  $m_{n+1}$ , it makes sense to use  
 214 the retrospective ratio which also involves the next model estimate  $m_{n+1}$  and not  
 215 the prospective ratio which involves the current model estimate  $m_n$ . Using the ret-  
 216 rospective ratio is however slightly more expensive because the next (approximate)  
 217 Hessian operator in the current direction must be computed in addition. Moreover

218 the accuracy of the retrospective prediction might be good in the direction  $-p_n$  while  
 219 still being bad in the direction  $p_{n+1}$  and inversely. There is also no safeguards for  
 220 large value of the ratios, which means that when the model is not accurate but the  
 221 predicted decrease underestimates the true decrease, the radius can still be increased.  
 222 Three sets of values for the threshold  $\rho_1$  and the rates  $c_0/c_1$  have been tested. The  
 223 acceptance threshold  $\rho_0$  is always tiny such that steps are often accepted, similarly to  
 224 the line search algorithm.

- 225 (A)  $\rho_0 = 10^{-4}$ ,  $\rho_1 = 0.25$  and  $c_0 = 0.20$ ,  $c_1 = 5$ .  
 226 (B)  $\rho_0 = 10^{-4}$ ,  $\rho_1 = 0.75$  and  $c_0 = 0.25$ ,  $c_1 = 2$ .  
 227 (C)  $\rho_0 = 10^{-4}$ ,  $\rho_1 = 0.90$  and  $c_0 = 0.50$ ,  $c_1 = 2$ .

228 The first one (A) is very similar to what was originally proposed in [10]. The other  
 229 two (B,C) are more cautious because they modify the radius more rarely and when  
 230 they do, it increases by a smaller factor. Note that the second one (B) is also close  
 231 to what is proposed in [23].

232 **2.2. Inner product.** The choice of the inner product plays a central role in the  
 233 inversion because it defines through the norm how directions length are measured but  
 234 also because it defines both gradients and Hessian operators. Indeed the equivalence  
 235 between both expansions (2.1) and (2.2) is granted by the defining property of the  
 236 gradient and the Hessian operator in terms of directional derivatives

237 (2.13)  $\langle j', \delta m_1 \rangle_M := \{D_m J\}(\delta m_1) \quad \forall \delta m_1,$

238 (2.14)  $\langle H\delta m_2, \delta m_1 \rangle_M := \{D_{mm}^2 J\}(\delta m_1, \delta m_2) \quad \forall \delta m_1, \delta m_2.$

240 This link between directional derivatives and kernels is actually a straightforward  
 241 application of the Fréchet-Riesz representation theorem [15].

242 The model parameter space is a function space defined on some region  $\Omega$  and  
 243 conventionally, the inner product is chosen as the  $L_2(\Omega)$  inner product

244 (2.15)  $\langle m_2, m_1 \rangle_M = \langle m_2, m_1 \rangle := \int_{\Omega} m_1(\mathbf{x})m_2(\mathbf{x}) d\Omega.$

245 This straightforward choice leads to the conventional gradient  $j'_{L_2}$  and the conven-  
 246 tional Hessian operator  $H_{L_2}$ , that can both be computed efficiently using the adjoint  
 247 state method [1, 13, 28]. As an illustration, a conventional gradient is represented  
 248 in Fig. 1b. It is actually the first gradient computed during the acoustic imaging  
 249 of the Marmousi model. As can be seen, shallow contributions have much greater  
 250 amplitudes than deeper parts. This actually reflects the bad scaling properties of this  
 251 inner product and motivates the use of a spatially weighted inner product

252 (2.16)  $\langle m_2, m_1 \rangle_M := \langle m_2 \sqrt{w}, \sqrt{w} m_1 \rangle,$

253 with an appropriate spatially dependant weight  $w$ . Insights on how to design  $w$  can  
 254 be gained by relating the conventional and the weighted gradients. Indeed, both are  
 255 defined by (2.13) then by transitivity of the equality

256 (2.17)  $\langle j'_{L_2}, \delta m_1 \rangle = \langle j' \sqrt{w}, \sqrt{w} \delta m_1 \rangle \quad \forall \delta m_1 \quad \text{such that} \quad j' = w^{-1} j'_{L_2}.$

257 The same reasoning can be applied to both Hessian operators ( $H = w^{-1} H_{L_2}$ ). Choos-  
 258 ing this weight close to the Hessian operator then makes the gradient closer to the  
 259 pure Newton direction and the Hessian operator closer to the identity. In other words,  
 260 the Newton system (2.3) is better conditioned and iterative solvers are therefore ex-  
 261 pected to converge faster. We choose here to take this weight as the diagonal part  
 262 of the Gauss-Newton Hessian ( $w = \text{diag}(H_{GN})$ ) because it can be computed semi-  
 263 analytically for a given model at no extra computational cost under certain circum-  
 264 stances [25]. A weight that has the same units than the Hessian also has the advantage

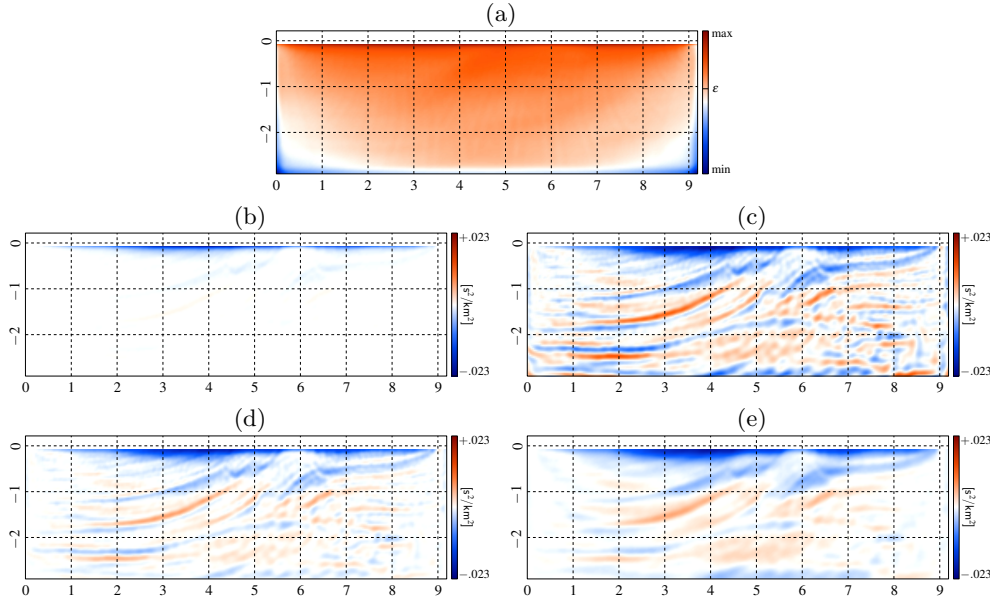


Fig. 1: Diagonal part of the Gauss-Newton Hessian (a). Conventional gradient (b). Weighted gradient (c). Weighted and thresholded gradient (d). Weighted and smoothed gradient (e). The stabilization parameter  $\epsilon$  is given graphically in the top figure and a smoothing length  $2\pi l_c = 0.250$  [km] is chosen.

265 that the corresponding weighted gradient has the same units than the model parameters. Model parameters, weighted gradients and weighted Hessian vector products  
 266 therefore all have the units of model parameters and the coefficients between them,  
 267 for example the length  $\gamma$  and  $\mu$  involved respectively in line search and trust region  
 268 techniques, are then always dimensionless and thus easier to interpret. The weights  
 269 and the corresponding weighted gradient are given in Fig. 1a and 1c respectively. As  
 270 expected, the weighted inner product compensates for the geometrical spreading and  
 271 restores balance between shallow and deep contributions. It is however dangerous  
 272 to use this weight alone because it can be very close to zero in poorly illuminated  
 273 zones as for example in the corners of the model. In these regions, the weighted inner  
 274 product is insensitive and consequently the preconditioner is unstable. The simplest  
 275 stabilization strategy consists in the introduction of a threshold  $\epsilon$  in the weights  
 276

$$277 \quad (2.18) \quad \langle m_2, m_1 \rangle_M := \langle m_2 \sqrt{w}, \sqrt{w} m_1 \rangle + \epsilon \langle m_2, m_1 \rangle.$$

278 The corresponding preconditioning effect is to keep silent regions where the weight is  
 279 much smaller than the threshold. Another strategy is to use an inner product with  
 280 the following stabilization term

$$281 \quad (2.19) \quad \langle m_2, m_1 \rangle_M := \langle m_2 \sqrt{w}, \sqrt{w} m_1 \rangle + \epsilon l_c^2 \langle \nabla m_2, \nabla m_1 \rangle$$

282 where  $l_c$  is a characteristic length. This second term, related to spatial derivation,  
 283 increases the norm of directions that are rapidly varying and also prevents the inner  
 284 product from being insensitive in regions where the diagonal Hessian is close to zero.  
 285 In regions where the diagonal Hessian is close to the threshold, then directions with  
 286 details smaller than the characteristic length  $l_c$  are penalized with respect to smoother  
 287 directions. This inner product is actually very similar to the one introduced in [41],  
 288 except that the Gauss-Newton diagonal Hessian weight is used in addition. As far  
 289 as preconditioning is concerned, this inner product can be reformulated through an

290 integration by parts as

$$291 \quad (2.20) \quad \langle m_2, m_1 \rangle_M := \langle w m_2, m_1 \rangle + \epsilon l_c^2 \langle \Delta m_2, m_1 \rangle.$$

292 Then as previously, conventional and preconditioned gradients are linked

$$293 \quad (2.21) \quad \langle j', \delta m_1 \rangle_M = \langle j'_{L_2}, \delta m_1 \rangle \quad \forall \delta m_1$$

$$294 \quad (2.22) \quad \langle (w + \epsilon l_c^2 \Delta) j', \delta m_1 \rangle = \langle j'_{L_2}, \delta m_1 \rangle \quad \forall \delta m_1 \quad \Leftrightarrow \quad j' = (w + \epsilon l_c^2 \Delta)^{-1} j'_{L_2}.$$

296 From the point of view of preconditioning, this inner product generates a rescaling  
 297 thanks to the Gauss-Newton diagonal Hessian weight and a Laplacian filtering, whose  
 298 smoothing length equals  $2\pi l_c$  where the diagonal Hessian equals the threshold. The  
 299 effect of these inner products is illustrated in Fig. 1d and 1e. In addition of stabilizing  
 300 the weights, [41] have shown that a filtering inner product can help the convergence  
 301 of full waveform inversion by mitigating its non linearity.

302 In general, any inner product that can be related to the conventional inner prod-  
 303 uct (2.15) through some preconditioner  $P$  yields a preconditioned gradient and a  
 304 preconditioned Hessian operator

$$305 \quad (2.23) \quad \langle m_2, m_1 \rangle_M = \langle P m_2, m_1 \rangle \quad \Rightarrow \quad j' = P^{-1} j'_{L_2} \text{ and } H = P^{-1} H_{L_2}.$$

306 Changing the inner product is formally equivalent to preconditioning both the gradient  
 307 and the Hessian operator. We choose to introduce preconditioning through a change  
 308 in the inner product rather than through the application of an operator because it  
 309 appears more elegant and rigorous to us. Moreover, this approach has the pedagogical  
 310 advantage to include preconditioning inside the inner product choice and thus it does  
 311 not need to appear explicitly in the description of the optimization algorithms. In  
 312 terms of practical implementation, it implies that the optimization routines must  
 313 not be rewritten, only the subroutine which computes the inner product have to be  
 314 modified, hence providing a lot of flexibility. Basically, a different choice for the inner  
 315 product does not modify the pure Newton direction because the same preconditioner  
 316 is applied to both sides of the Newton system (2.3), but does modify the subspace  
 317 constructed by the conjugate gradient method and does modify norms which are  
 318 involved in any stopping criterion. A good choice can thus lead to better approximate  
 319 directions and better truncation rules.

320 **2.3. Steepest descent.** The steepest descent is actually the simplest local opti-  
 321 mization algorithm. It consists in taking the search direction as the opposite gradient.  
 322 This is the best direction at first order ( $\tilde{H} = 0$ ) but it can also be seen as a quasi-  
 323 Newton step where the approximate Hessian operator is the identity operator ( $\tilde{H} = I$ ).  
 324 In practice however, this approximation is very crude because the Hessian operator  
 325 is far from the identity operator, even after preconditioning. The downside of this  
 326 simple method is its linear convergence rate. This slow convergence speed is one of  
 327 the main motivation for the investigation of higher order algorithms.

328 **2.3.1. Line search globalization.** No length information can be captured  
 329 from the approximate Hessian operator in this case, because it is simply the iden-  
 330 tity operator ( $\tilde{H} = I$ ). The first trial step length is then chosen based on the  
 331 history of the outer iterations to save as many step length trials as possible *e.g.*  
 332  $\gamma = 2(J(m_n) - J(m_{n-1})) / \{D_m J\}(-j')$  [23].

333 **2.3.2. Trust region globalization.** Trust-region methods are barely used with  
 334 steepest descent. Mostly because the linear misfit prediction

$$335 \quad (2.24) \quad J^{\text{pred}}(m; p) := J(m) + \langle j'(m), p \rangle_M$$

336 is not accurate enough. Moreover the solution to the trust-region sub-problem (2.6) is  
 337 trivially  $p = -\mu j'$  and is always on the boundary, because of the absence of a second



338 order term. An upper bound on the relative size of the trust region ( $\mu$ ) is then added  
 339 to compensate the fact that the trust-region algorithm will never keep it constant.  
 340 This bound is set to  $\mu_{\max} = 4, 4, 5$  for parameter sets A, B, C respectively.

341 **2.4. Limited memory BFGS method.** Quasi-Newton methods are expected  
 342 to provide a huge improvement over the steepest descent and an attractive alternative  
 343 to Newton methods because they do not involve any expensive Hessian vector prod-  
 344 uct. In place of the exact Hessian, an approximation  $\bar{H} = B$  is used instead. This  
 345 approximation is built only with the successive gradients and model parameters of  
 346 each iteration. Moreover, since expensive Hessian vector product are avoided, quasi-  
 347 Newton methods are sometimes more efficient than Newton methods. The Broyden-  
 348 Fletcher-Goldfarb-Shanno algorithm, abbreviated BFGS, is maybe the most widely  
 349 used quasi-Newton method. This method constructs a symmetric and positive defi-  
 350 nite approximation of the Hessian operator based on all the previous gradients and  
 351 model parameters. This approximation  $B_{n+1}$  is chosen such that it verifies the secant  
 352 equation

$$353 \quad (2.25) \quad B_{n+1} \Delta m_n = \Delta j'_n \quad \text{with } \Delta m_n = m_{n+1} - m_n \text{ and } \Delta j'_n = j'_{n+1} - j'_n$$

354 while being close to the previous approximation  $B_n$  and positive definite. Note that  
 355 imposing the positive definiteness of this approximation also imposes that the update  
 356 direction must satisfy the (BFGS) curvature condition  $\langle \Delta m_n, \Delta j'_n \rangle_M > 0$ . One of the  
 357 biggest advantage of the BFGS algorithm is that it is possible to directly build the  
 358 approximate inverse Hessian operator  $B_n^{-1}$  from the memorized gradients and model  
 359 parameters. However building explicitly this inverse operator in the context of large-  
 360 scale optimization is still prohibitively expensive, as well as storing in memory all  
 361 the previous gradients and model parameters. For these reasons, a limited memory  
 362 version of the algorithm has been derived. Instead of memorizing all the previous  
 363 iterates, it only requires the  $l$  last iterates and above all, it comes with a two-loop re-  
 364 cursive procedure to compute the application of the inverse operator on any direction.  
 365 The approximate Newton direction associated with the  $l$ -BFGS operator is therefore  
 366 straightforward to compute. This two-loop recursive  $l$ -BFGS algorithm is given in  
 367 Algorithm 2.2 [23].

**Algorithm 2.2**Inverse  $l$ -BFGS operator application

---

**Require:**  $q, \Delta m_k, \Delta j'_k, \forall k \in [n-l, n-1]$   
**for**  $k = n-1$  **down to**  $k = n-l$  **do**  
 $\alpha_k = \langle \Delta m_k, q \rangle_M / \langle \Delta j'_k, \Delta m_k \rangle_M$   
 $q = q - \alpha_k \Delta j'_k$   
**end for**  
 $\xi = \langle \Delta m_{n-1}, \Delta j'_{n-1} \rangle_M / \langle \Delta j'_{n-1}, \Delta j'_{n-1} \rangle_M$   
 $r = \xi q$   
**for**  $k = n-l$  **up to**  $k = n-1$  **do**  
 $\beta_k = \langle \Delta j'_k, r \rangle_M / \langle \Delta j'_k, \Delta m_k \rangle_M$   
 $r = r + (\alpha_k - \beta_k) \Delta m_k$   
**end for**  
**return**  $r (= B_n^{-1} q)$

---

**Algorithm 2.3**Direct  $l$ -BFGS operator application

---

**Require:**  $q, \Delta m_k, \Delta j'_k, \forall k \in [n-l, n-1]$   
**for**  $k = n-l$  **up to**  $k = n-1$  **do**  
 $b_k = \Delta j'_k / \sqrt{\langle \Delta j'_k, \Delta m_k \rangle_M}$   
 $a_k = B_n^0 \Delta m_k$   
**for**  $i = n-l$  **up to**  $i = k-1$  **do**  
 $a_k = a_k + \langle b_i, \Delta m_k \rangle b_i - \langle a_i, \Delta m_k \rangle a_i$   
**end for**  
 $a_k = a_k / \sqrt{\langle \Delta m_k, a_k \rangle_M}$   
**end for**  
 $r = B_n^0 q$   
**for**  $k = n-l$  **up to**  $k = n-1$  **do**  
 $r = r + b_k \langle b_k, q \rangle_M - a_k \langle a_k, q \rangle_M$   
**end for**  
**return**  $r (= B_n q)$

---

369 It is important to highlight here that this method also benefits from the modification  
 370 of the inner product. Indeed the building blocks of this approximate Hessian operator  
 371 are the successive gradients, which are preconditioned through the inner product. By  
 372 measuring gradient variations, this method constructs a representation of the misfit  
 373 which is good enough to produce super-linear convergence, a great improvement over

374 the steepest descent, at no extra cost. This approximation is however positive definite  
 375 while the exact Hessian might be indefinite, especially during the early iteration of  
 376 the inversion. In such cases, this quasi-Newton method may fail to converge while  
 377 Newton methods may not [20].

378 **2.4.1. Line search globalization.** The unit step length  $\gamma = 1$  is always tried  
 379 first because the length information should be captured by the inverse approximate  
 380 Hessian. Importantly, it can be showed that the (BFGS) curvature condition is always  
 381 satisfied if the strong Wolfe conditions (2.4) and (2.5) are enforced [23]. Therefore  
 382 the  $l$ -BFGS algorithm combined with a line search will always construct a positive  
 383 definite approximate Hessian operator  $B$ .

384 **2.4.2. Trust region globalization.** Finding the exact solution to the trust-  
 385 region sub-problem (2.6) with the  $l$ -BFGS predicted misfit

$$386 \quad (2.26) \quad J^{\text{pred}}(m; p) := J(m) + \langle j'(m), p \rangle_M + 0.5 \langle Bp, p \rangle_M$$

387 is difficult for a general trust region radius. However when this radius is large enough,  
 388 in particular larger than the unconstrained solution  $p^u := -B^{-1}j'$ , then it is ac-  
 389 tually also the exact constrained solution. On the other hand, when the radius is  
 390 small enough, the quadratic term in the misfit prediction is negligible and the sub-  
 391 problem is equivalent to the steepest descent, which indicates to follow the gradient  
 392 up to the boundary. Based on these solutions for extreme value of the radius, the  
 393 exact solution to the sub-problem (2.6) might be substituted by an interpolation be-  
 394 tween these two solutions. Namely, the gradient is followed each time the minimum  
 395 of the misfit prediction along the gradient, *i.e.* the Cauchy point  $p^c = -\alpha j'$  (with  
 396  $\alpha = \langle j', j' \rangle_M / \langle B j', j' \rangle_M$ ), is outside the radius. Then for intermediate radii, which  
 397 contains this Cauchy point but not the unconstrained solution, an interpolation be-  
 398 tween both is done

$$399 \quad (2.27) \quad p(\Delta) = p^c + \tau^* (p^u - p^c) \quad \text{with } 0 < \tau^* < 1 \text{ such that } \|p\|_M = \Delta.$$

400 Finally for large radii, the unconstrained solution is accepted. In summary

$$401 \quad (2.28) \quad p(\Delta) = \begin{cases} p^u & \text{when } \|p^u\|_M \leq \Delta, \\ -\mu j' & \text{when } \|p^c\|_M \geq \Delta, \\ p^c + \tau^* (p^u - p^c) & \text{when } \|p^c\|_M \leq \Delta \leq \|p^u\|_M. \end{cases}$$

402 The approximate solution (2.28) to the trust-region sub-problem (2.6) is called the  
 403 dogleg method [23].

404 A huge difference with the line search implementation of the  $l$ -BFGS algorithm  
 405 is that now the direct application of the approximate Hessian operator  $B$  on some  
 406 directions must be computed. Unfortunately there is no equivalent to Algorithm 2.2  
 407 for the direct  $l$ -BFGS operator and its application must then be computed from its  
 408 recursive definition

$$409 \quad (2.29) \quad B_n q = B_n^0 q + \sum_{k=n-l}^{n-1} b_k \langle b_k, q \rangle_M - a_k \langle a_k, q \rangle_M$$

$$410 \quad (2.30) \quad \text{with } a_k = \frac{B_k \Delta m_k}{\sqrt{\langle B_k \Delta m_k, \Delta m_k \rangle_M}} \quad \text{and} \quad b_k = \frac{\Delta j'_k}{\sqrt{\langle \Delta j'_k, \Delta s_k \rangle_M}}.$$

412 It is important to highlight that the sequence of directions  $a_k$  could not be memo-  
 413 rized because at each iterations the oldest information is discarded, which modifies  
 414 the whole  $a_k$  sequence. A complete procedure to compute the application of the direct

415  $l$ -BFGS operator is given in Algorithm 2.3. Faster but more sophisticated procedure  
 416 do exist [23]. However manipulations in the model parameter space are computationally  
 417 negligible with respect to wave propagation problems hence the speedup would  
 418 also be negligible. Thanks to this procedure the prospective and retrospective predicted  
 419 decrease (2.10) and (2.12) can be evaluated. Interestingly, the prospective  
 420 decrease is evaluated with the current Hessian approximation  $B_n$  while the retrospective  
 421 decrease is evaluated with the next Hessian approximation  $B_{n+1}$ . The retro-  
 422 spective ratio is therefore expected to be more often close to one because this next  
 423 Hessian approximation  $B_{n+1}$  is specifically constructed from the update direction  
 424  $p_n = \Delta m_n = m_{n+1} - m_n$ .

425 **2.5. Newton methods.** In contrast with quasi-Newton methods, Newton meth-  
 426 ods use the Hessian operator explicitly, as they try to solve the Newton system (2.3).  
 427 The interest of these method lies in their independence on the convergence history  
 428 and in their quadratic convergence rate in the vicinity of the minimum. Far from this  
 429 minimum, the Hessian operator might however be indefinite, which complicates the  
 430 solution procedure for the Newton system. For that reason, it is frequent to make the  
 431 Gauss-Newton approximation ( $\tilde{H} = H_{GN}$ ), which consist in keeping only the positive  
 432 definite part of the Hessian operator. The downside of this approximation is then  
 433 that the second order representation (2.2) of the misfit is less accurate, especially if  
 434 the negative definite part of the Hessian is not negligible, which might prevent the  
 435 method from converging. In this section, we present inexact Newton methods based  
 436 on a line search procedure or a trust region method. Both are valid for the full Hessian  
 437 and for its Gauss-Newton approximation.

---

**Algorithm 2.4** Conventional  
conjugate gradient algorithm

---

$p_0 = 0, r_0 = j', q_0 = -j'$   
**if**  $\langle H j', j' \rangle_M \leq 0$  **then return**  $-j'$   
**loop**  
**if**  $\langle H q_k, q_k \rangle_M \leq 0$  **then return**  $p_k$

$$\alpha_k = \langle r_k, r_k \rangle_M / \langle H q_k, q_k \rangle_M$$

$p_{k+1} = p_k + \alpha_k q_k$  and  $r_{k+1} = r_k + \alpha_k H q_k$   
**if**  $\|r_{k+1}\|_M < \eta \|j'\|_M$  **then return**  $p_{k+1}$   
 $\beta_{k+1} = \|r_{k+1}\|_M^2 / \|r_k\|_M^2$   
 $q_{k+1} = -r_{k+1} + \beta_{k+1} q_k$

**end loop**

---



---

**Algorithm 2.5** Steihaug  
conjugate gradient algorithm

---

$p_0 = 0, r_0 = j', q_0 = -j'$   
**loop**  
**if**  $\langle H q_k, q_k \rangle_M \leq 0$  **then**  
 $\tau^* = \tau > 0 \mid \|p_k + \tau q_k\|_M = \Delta$   
**return**  $p_k + \tau^* q_k$   
**end if**

$\alpha_k = \langle r_k, r_k \rangle_M / \langle H q_k, q_k \rangle_M$   
**if**  $\|p_k + \alpha_k q_k\|_M \geq \Delta$  **then**  
 $\tau^* = \tau > 0 \mid \|p_k + \tau q_k\|_M = \Delta$   
**return**  $p_k + \tau^* q_k$   
**end if**

$p_{k+1} = p_k + \alpha_k q_k$  and  $r_{k+1} = r_k + \alpha_k H q_k$   
**if**  $\|r_{k+1}\|_M < \eta \|j'\|_M$  **then return**  $p_{k+1}$   
 $\beta_{k+1} = \|r_{k+1}\|_M^2 / \|r_k\|_M^2$   
 $q_{k+1} = -r_{k+1} + \beta_{k+1} q_k$

**end loop**

---

439 **2.5.1. Line search globalization.** Newton methods can be combined with a  
 440 line search procedure. In this case a direction  $p$  is first found by solving the Newton  
 441 system approximately with the conventional conjugate gradient method (Algorithm  
 442 2.4) [23]. This algorithm constructs iteratively the solution of a linear system without  
 443 requiring the explicit expression of the Hessian matrix but only its action in particular  
 444 directions. The iterative procedure is stopped when the residuals have decreased more  
 445 than some threshold, called the forcing sequence  $\eta$ , which is typically close to zero

446 (2.31)  $(\|r_k\|_M :=) \|H p_k + j'\|_M < \eta \|j'\|_M (= \eta \|r_0\|_M).$

447 Over-solving is here avoided through this forcing term  $\eta$ , which is not systematically  
 448 close to zero but which is instead chosen to reflect the accuracy of the second-order  
 449 expansion. Three possible choices for this sequence have been described and studied  
 450 by [9]. These three choices were then compared in the context of acoustic imaging in  
 451 [20], who advise to use the forcing sequence

$$452 \quad (2.32) \quad \eta_m = \frac{\|j'(m_n) - j'(m_{n-1}) - \gamma_{n-1}H(m_{n-1})p_{n-1}\|_M}{\|j'(m_{n-1})\|_M}.$$

453 If the accuracy of the local quadratic approximation is good then this forcing term  
 454 is close to zero and the Newton system is solved accurately. If not, then iterations  
 455 are truncated sooner. This forcing sequence plays a similar role than the prospective  
 456 ratio for trust-region method. It is however based on a (prospective) expansion of the  
 457 gradient while the prospective ratio is based on an expansion of the misfit. Additional  
 458 safeguards are also added to prevent this forcing term to decrease too fast or to increase  
 459 above  $\eta_0 = 0.9$ . Interestingly, directions of negative curvatures are never investigated,  
 460 except if it is the gradient. As previously an appropriate length  $\gamma$  is then given to  
 461 this direction  $p$  through a line search. The unit step length  $\gamma = 1$  is again tried first  
 462 because it is the best choice if the misfit were quadratic.

463 **2.5.2. Trust region globalization.** When the Newton method is associated  
 464 with a trust-region technique, the direction is found by minimizing the local quadratic  
 465 expansion of the misfit

$$466 \quad (2.33) \quad J^{\text{pred}}(p) := J(m) + \langle j', p \rangle_M + 0.5 \langle Hp, p \rangle_M$$

467 inside a sphere of radius  $\Delta$ . The constraint  $\|p\|_M \leq \Delta$  limits the size of the direction  
 468 and aims at preventing over-solving. This trust-region sub-problem can be solved ap-  
 469 proximately with the Steihaug conjugate gradient method (Algorithm 2.5) [33]. This  
 470 method actually exploits two properties of the conjugate gradient algorithm: succes-  
 471 sive approximate solutions always grow in norm ( $\|p_k\|_M < \|p_{k+1}\|_M$ ) while the misfit  
 472 prediction always decrease ( $J^{\text{pred}}(p_k) > J^{\text{pred}}(p_{k+1})$ ). The underlying idea of the  
 473 method is then to minimize the second order expansion of the misfit iteratively using  
 474 the conventional conjugate gradient algorithm until either convergence is achieved, ei-  
 475 ther the boundary is reached. Basically there are only two modifications compared  
 476 to Algorithm 2.4. First, the inner iterations are cropped to the trust region radius  
 477  $\Delta$  when the unconstrained solution increases beyond it. Second, when a direction of  
 478 negative curvature is encountered, it is followed up to the boundary of the trust region  
 479 and the algorithm is stopped. Interestingly these directions were never investigated in  
 480 the conventional version. The convergence criterion is unchanged but here the forcing  
 481 term is kept constant ( $\eta = 0.5$ ).

482 **3. Numerical investigations.** Numerical studies are performed in the context  
 483 of subsurface acoustic imaging in the frequency domain [29, 32]. In that particular  
 484 case, the misfit is conventionally chosen as the least-squares distance between some  
 485 acoustic pressure measurements  $d_{\omega er}$  (at some receiver  $r$ , for several excitation sources  
 486  $e$  and for different frequencies  $\omega$ ) and the corresponding computed acoustic pressures  
 487  $p_{\omega e}(\mathbf{x}_r)$ , obtained by solving the Helmholtz equation

$$488 \quad (3.1) \quad J(s^2) = 0.5 \sum_{\omega, e, r} |p_{\omega e}(\mathbf{x}_r; s^2) - d_{\omega er}|^2 \quad \text{with} \quad \Delta p + \omega^2 s^2 p = \delta(\mathbf{x} - \mathbf{x}_e).$$

489 It is here chosen that the subsurface model parameter is the slowness squared distri-  
 490 bution  $s^2$  [ $\text{s}^2/\text{km}^2$ ] (also called the sloth), as could be guessed from the expression of  
 491 the Helmholtz operator  $A_\omega(s^2) := \Delta + \omega^2 s^2$ . The slowness squared  $s^2$  is actually the  
 492 squared inverse of the velocity  $v$ . Several other parametrizations are also possible but

493 it has been shown that the slowness squared can yield a fast convergence and accurate  
 494 results [2, 5, 14, 27]. Implementation of any of the above described local optimization  
 495 algorithms requires an efficient procedure to compute the misfit and the gradient for a  
 496 given slowness squared distribution  $s^2$  and the action of the Hessian operator for any  
 497 given slowness squared perturbation  $\delta s^2$ . The well-known adjoint state method has  
 498 been developed for that specific purpose. It is summarized here below and detailed  
 499 in [1, 13, 28]. The two terms in gray should be removed under the Gauss-Newton  
 500 approximation.

- 501 1. Find the forward fields  $p_{\omega e}$  such that

$$502 \quad (3.2) \quad A_{\omega}(s^2)p_{\omega e} = \delta(\mathbf{x} - \mathbf{x}_e).$$

- 503 2. Find the adjoint fields  $p_{\omega e}^{\dagger}$  such that

$$504 \quad (3.3) \quad A_{\omega}(s^2)p_{\omega e}^{\dagger} = \sum_r (\bar{p}_{\omega e}(\mathbf{x}_r) - \bar{d}_{\omega e r})\delta(\mathbf{x} - \mathbf{x}_r).$$

- 505 3. Find the preconditioned gradient  $j'$  such that

$$506 \quad (3.4) \quad Pj' = - \sum_{\omega} \omega^2 \sum_e p_{\omega e}^{\dagger} \bar{p}_{\omega e}.$$

- 507 4. Find the perturbed forward fields  $\delta p_{\omega e}$  such that

$$508 \quad (3.5) \quad A_{\omega}(s^2)\delta p_{\omega e} = -\omega^2 \delta s^2 p_{\omega e}.$$

- 509 5. Find the perturbed adjoint fields  $\delta p_{\omega e}^{\dagger}$  such that

$$510 \quad (3.6) \quad A_{\omega,e}(s^2)\delta p_{\omega e}^{\dagger} = \sum_r \delta p_{\omega e}(\mathbf{x}_r)\delta(\mathbf{x} - \mathbf{x}_r) - \omega^2 \delta s^2 p_{\omega e}^{\dagger}.$$

- 511 6. Find the preconditioned Hessian operator  $H\delta s^2$  in the direction  $\delta s^2$  such that

$$512 \quad (3.7) \quad PH\delta s^2 = - \sum_{\omega} \omega^2 \sum_e (\delta p_{\omega e}^{\dagger} \bar{p}_{\omega e} + p_{\omega e}^{\dagger} \delta \bar{p}_{\omega e}).$$

513 Independently of any practical solver for these wave propagation problems, a misfit  
 514 evaluation only requires to perform step 1 and thus only requires to solve a single  
 515 wave propagation problem. A gradient evaluation requires steps 1 to 3, thus a single  
 516 supplementary wave propagation problem must be solved if the misfit has already been  
 517 computed. Similarly, steps 1 to 6 are necessary for the application of the (Gauss-  
 518 )Newton Hessian operator in a particular direction, thus again two supplementary  
 519 wave propagation problems if the gradient has already been computed for the same  
 520 model parameters.

521 Consequently the steepest descent and the  $l$ -BFGS directions require to solve  
 522 two wave problems while any Newton-based direction has an initial cost of four wave  
 523 propagation problems and each supplementary conjugate gradient iteration requires  
 524 two more wave problems. To the price of the directions must be added the cost  
 525 of the line search or the trust-region methods. Line search typically accepts a step  
 526 length if it verifies sufficient conditions (2.4) and (2.5) which involves the misfit  
 527 and its gradient. Thus it requires one or two additional wave problems each time a trial  
 528 step length is rejected. Prospective trust-region has no additional cost because the  
 529 evaluation of the trust region only depends on quantities already computed. At the  
 530 opposite, retrospective (Gauss-)Newton trust-region requires the application of the  
 531 Hessian operator on the preceding direction and thus needs to solve two additional  
 532 wave propagation problems. The table here below summarizes this accounting.

	Base	CG	LS	TR
SD	2	-	$2N_{LS}$	-
$l$ -BFGS	2	-	$2N_{LS}$	-
LS-NCG	2	-	$2N_{LS}$	-
TR-NCG (P)	2	$2N_{CG}$	-	0
TR-NCG (R)	2	$2N_{CG}$	-	2

545

546 It is interesting to highlight here that the first inner iteration of any conjugate gradient  
 547 Newton method is simply the steepest descent but it is twice more expensive because  
 548 the curvature is computed. Subsequent inner iterations must therefore provide large  
 549 decrease of the misfit to compensate this high entry cost. This phenomenon is even  
 550 worse with the retrospective trust region algorithm because there is a systematical  
 551 additional cost to update the trust region radius.

552 In this work, solutions to partial differential equations (3.2) to (3.7) are obtained  
 553 numerically with the finite element method. In what follows, we specify the exact  
 554 numerical procedure in that context. Note however that the analysis would nearly  
 555 be identical with finite differences. Finite element discretization assembles operators  
 556 into matrices and source terms into vectors. Wave propagation problems (3.2), (3.3),  
 557 (3.5) and (3.6) therefore transform into a linear system whose left-hand-side matrix  
 558  $A$  is always the same for a given frequency while the right-hand-side source  $b$   
 559 is different for any field type, frequency and excitation index. The solution of this  
 560 system is obtained by first computing its lower-upper factorization then by performing  
 561 an upward-backward substitution for each right-hand-side source

$$562 \quad (3.8) \quad Ap = b \quad \Leftrightarrow \quad A = LU, \quad Lq = b \text{ and } Up = q.$$

563 Huge computational reduction is therefore obtained because only one matrix per fre-  
 564 quency is assembled and factorized. The computation of any wave field then requires  
 565 the assembly and the upward-backward substitution of a vector per excitation source,  
 566 but no more matrix factorization. The numerical equivalence of the preceding six  
 567 steps procedure is given here below.

- |     |    |                                        |                         |
|-----|----|----------------------------------------|-------------------------|
| 568 | 1. | • Factorize wave propagation operators | $(n_\omega)$            |
| 569 |    | • Substitute forward sources           | $(n_\omega \times n_e)$ |
| 570 | 2. | • Substitute adjoint sources           | $(n_\omega \times n_e)$ |
| 571 | 3. | • Factorize the preconditioner         | (1)                     |
| 572 |    | • Substitute the conventional gradient | (1)                     |
| 573 | 4. | • Substitute perturbed forward sources | $(n_\omega \times n_e)$ |
| 574 | 5. | • Substitute perturbed adjoint sources | $(n_\omega \times n_e)$ |
| 575 | 6. | • Substitute the conventional Hessian  | (1)                     |

576 It is interesting to highlight that model problems (steps 3 and 6) are negligible with  
 577 respect to wave problems. Indeed while wave problems involve a matrix per frequency  
 578 and a vector per excitation source, model problems only involve a single matrix (*i.e* the  
 579 preconditioner) and a single source vector (*i.e* the conventional gradient or Hessian).  
 580 Moreover the model discretization is usually coarser than the wave field discretization.  
 581 Consequently not considering these model problems when quantifying the computa-  
 582 tional complexity is not dramatic. It should however be highlighted that forward  
 583 problems are more expensive than the corresponding adjoint problem, because the  
 584 matrix factorization is reused. Moreover the perturbed forward problem and the per-  
 585 turbed adjoint problem are slightly heavier than the adjoint problem, because both  
 586 their sources are dense, at the opposite of forward and adjoint sources, which are  
 587 sparse. Nevertheless we weight equally all of these four problems when quantifying  
 588 the computational complexity.

589 In the next two sections, two synthetic numerical case studies are investigated.  
 590 The first one is based on the widely used Marmousi benchmark [34] while the sec-  
 591 ond one, replicated from [20], is inspired from a near-surface imaging of close concrete  
 592 structures and features important multiple scattering. Multiple scattering is responsi-  
 593 ble for the indefiniteness of the Hessian operator, which, as mentioned in the previous  
 594 part, is challenging for optimization algorithms. This second example is thus chosen  
 595 to emphasize which optimization methods are able to overcome such difficulties. For  
 596 both case studies, the influence of the inner product choice on the convergence speed  
 597 and the quality of the inverted model is studied first. Once the inner product is cho-  
 598 sen, prospective and retrospective trust-region methods with different parameter sets

599 are compared and the best option is selected. Advantages and drawbacks of trust-  
 600 region methods in the context of full waveform inversion are then finally discussed  
 601 based on a comparison with the corresponding line search methods. In the remainder  
 602 of this section, data misfit are normalized such that the misfit corresponding to the  
 603 initial model is one and computational complexity is measured in numbers of forward  
 604 problems solved, as explained above.

605 **3.1. Case study 1.** Numerical inversions are performed on the 2D Marmousi  
 606 model (Fig. 2a) [34] in the frequency domain. Three frequencies (4, 6 and 8 [Hz]) are  
 607 inverted simultaneously. The surface acquisition system is composed of 122 equally  
 608 spaced (72 [m]) excitation sources and 243 equally spaced (36 [m]) receivers. Outer  
 609 iterations are stopped when satisfying the convergence criterion  $J(s^2)/J(s_{\text{init}}^2) < 10^{-3}$ .  
 610 A smoothed version of the exact Marmousi model is used as an initial guess (Fig. 2b).  
 611 This initial model is computed with a Laplacian filter  $s_{\text{init}}^2 = (1 + (l_c/2\pi)^2 \Delta)^{-1} s_{\text{exact}}^2$   
 612 with  $l_c = 2$  [km]. Slowness squared fields and pressure fields at the three frequencies  
 613 are discretized on a square grid (36 [m]) by hierarchical finite elements, respectively  
 614 of order 1 and of order 2, 3, 4. A water layer (216 [m]) is also added at the top of the  
 615 model but it is kept constant during the inversion. The model is spatially truncated  
 616 by Sommerfeld boundary conditions [31]. Recorded data are generated synthetically  
 617 using the same hierarchical finite elements setting than for the inversion. An inversion  
 618 result, *i.e.* an estimated squared slowness, is shown in (Fig. 2c). From a relatively  
 619 low resolution initial guess, full waveform inversion indeed provides a high-resolution  
 620 estimation of the exact model. Images obtained with the other methods do not differ  
 621 significantly.

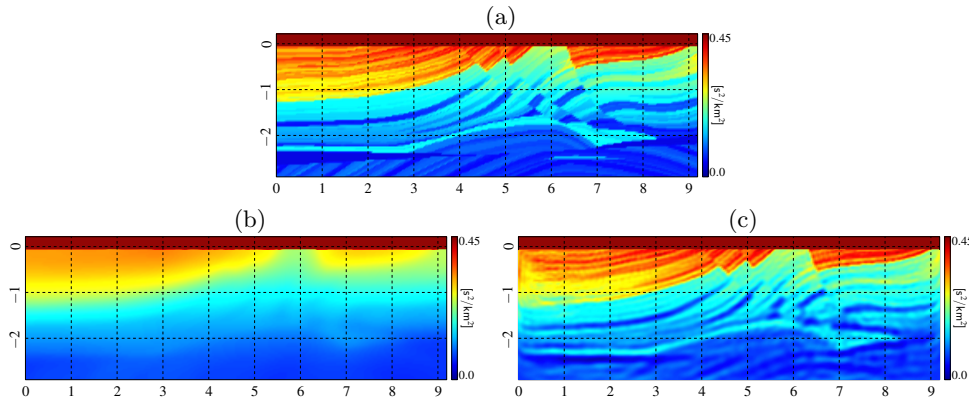


Fig. 2: Marmousi model (a), initial guess (b) and inversion results using a line search  $l$ -BFGS algorithm with a weighted and thresholded inner product (c).

622 **3.1.1. Inner product.** As explained earlier, the inner product has an influence  
 623 on both the gradient and the Hessian. Its choice is therefore expected to influence  
 624 the convergence speed but also the particular minimizer that is reached. To illustrate  
 625 both these effects, the line search  $l$ -BFGS algorithm has been applied with the four  
 626 different inner products introduced in this work, *i.e.* the conventional inner prod-  
 627 uct (2.15), the weighted inner product (2.16) and its regularized variants (2.18) or  
 628 (2.19). Corresponding convergence curves and error maps are given in Fig. 3 and  
 629 4 respectively. Both these figures are also summarized in Table 1. As can be seen  
 630 from these figures, all these weighted inner product increase the convergence speed with  
 631 respect to the conventional, *i.e.* unweighted, one. However the minimizer obtained  
 632 with the weighted inner product alone is further away from the exact solution, in  
 633 particular in the right corner of the model. Avoiding such artifacts is precisely one

634 of the reasons for the introduction of regularized inner products, as they dampen  
 635 the contributions in these poorly illuminated regions. Both the thresholding and the  
 636 smoothing strategy perform similarly in reducing the error back to the same level than  
 637 the unweighted solution but the thresholding strategy converges faster. It is thus kept  
 638 for the sequel of this case study. The advantages of the smoothing inner product will  
 639 be highlighted during the second case study. In the next three subsections, the be-  
 640 haviour of the steepest descent method, the  $l$ -BFGS method, the full Newton and the  
 641 Gauss-Newton methods is analysed. Convergence curves and interesting statistics for  
 all these methods are given in Fig. 5 and Table 2 respectively.

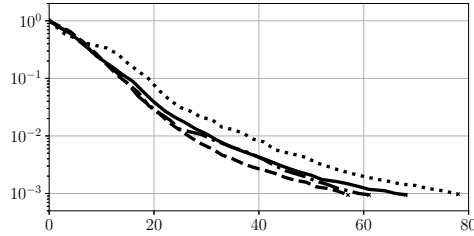


Fig. 3: Data misfit as a function of the computational complexity for the line search  $l$ -BFGS algorithm with a conventional ( $\bullet\bullet$ ), only weighted ( $-\bullet$ ), weighted and stabilized ( $--$ ) or weighted and smoothed ( $-$ ) inner product.

	Wave sol. (tot)	Error rms ( $[s^2/km^2]$ )
Conventional	78	0.0174
Weighted only	61	0.0202
and stabilized	57	0.0174
and smoothed	68	0.0173

Table 1: Computational complexity and root-mean squared error for the line search  $l$ -BFGS algorithm with different inner products.

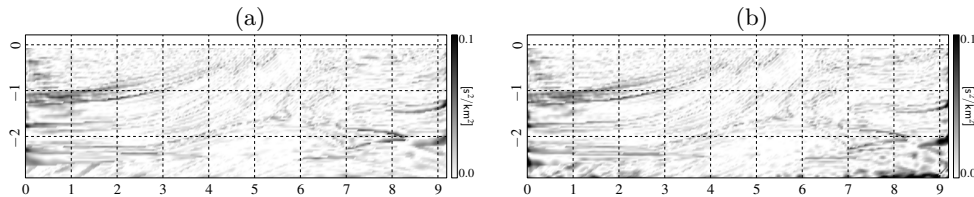


Fig. 4: Final inversion error for the line search  $l$ -BFGS algorithm with a conventional (a) or a weighted (b) inner product. Inversion errors for both regularized inner products are not shown because these are very similar to those obtained with the conventional inner product.

643 **3.1.2. Steepest descent.** There is no dramatic improvement when using one  
 644 or another direction scaling method, because actually the direction itself is bad. Nev-  
 645 ertheless, it appears that methods which reject less frequently the proposed update  
 646 direction are faster, *i.e* the prospective trust-region method with the more cautious  
 647 parameters sets (B and C) and the line search method. Retrospective radius update  
 648 does not speed up convergence. Actually we observed that the retrospective predicted  
 649 decrease (2.12) sometimes largely underestimates the actual decrease, illustrating that  
 650 the retrospective misfit prediction is very not accurate, but still producing an increase  
 651 of the trust region radius. Finally, among the three best methods, the slope is slightly  
 652 steeper for the two trust-region methods, probably because they systematically try to  
 653 increase the length given to the gradient direction.



		Wave sol. (tot)	Outer it. (tot)	Inner it. (avg)	Rejected (%)	Constrained (%)	Negative curv. (%)
SD	LS	244	111	10	-	-	-
	TR-P (A)	396	198	-	40	100	-
	TR-P (B)	280	140	-	6	100	-
	TR-P (C)	264	132	-	5	100	-
	TR-R (A)	328	164	-	20	100	-
	TR-R (B)	354	177	-	20	100	-
	TR-R (C)	330	165	-	25	100	-
LB	LS	57	27	-	7	-	-
	TR-P (A)	58	29	-	3	10	-
	TR-P (B)	58	29	-	3	34	-
	TR-P (C)	64	32	-	13	50	-
	TR-R (A)	58	29	-	3	10	-
	TR-R (B)	56	28	-	0	11	-
	TR-R (C)	56	28	-	0	11	-
FN	LS	139	17	2.9	12	-	29
	TR-P (A)	178	22	3.0	32	77	0
	TR-P (B)	106	13	3.1	0	69	0
	TR-P (C)	106	16	2.3	0	75	0
	TR-R (A)	144	14	3.1	14	64	0
	TR-R (B)	128	14	2.6	0	79	0
	TR-R (C)	142	17	2.2	0	82	0
GN	LS	124	15	3.13	0	-	-
	TR-P (A)	130	11	4.9	0	10	-
	TR-P (B)	98	10	3.9	0	30	-
	TR-P (C)	98	10	3.9	0	30	-
	TR-R (A)	152	11	4.9	0	10	-
	TR-R (B)	132	14	2.7	0	79	-
	TR-R (C)	184	24	1.8	0	83	-

Table 2: Statistics related to the implementation of the steepest descent (SD), the  $l$ -BFGS (LB), the full Newton (FN) method and the Gauss-Newton (GN) methods combined with a line search (LS) or a trust region (TR) with a prospective (P) or retrospective (R) radius update with different parameter sets (A,B,C).

654 **3.1.3. Limited memory BFGS method.** There is hardly no difference between  
655 all the methods combined with the  $l$ -BFGS algorithm. We observed that the  
656 line search method only rejects the unit step length  $\gamma = 1$  for the first two iterations.  
657 Similarly, we observed that the retrospective ratio is always very close to one, such  
658 that the trust region radius for retrospective methods quickly becomes large and thus  
659 the pure  $l$ -BFGS direction is always accepted after the first few iterations. An algo-  
660 rithm that unconditionally follows the pure  $l$ -BFGS direction would therefore already  
661 be very good and neither a line search nor a trust-region method can actually dras-  
662 tically improve it, as far as convergence speed is concerned. Nevertheless the more  
663 cautious prospective trust-region methods (B,C) also converge fast, which shows that,  
664 on the other hand, constraining the size of the update directions does not slow down  
665 the inversion.

666 **3.1.4. Newton methods.** As far as trust-region methods are concerned, it first  
667 clearly appears that the retrospective radius update is not worth its computational  
668 cost. Indeed it does not require less wave solutions than the best prospective ones,  
669 even if the computation cost of the retrospective predicted decrease is withdrawn (two  
670 wave solutions per outer iterations). Retrospective radius update has been introduced

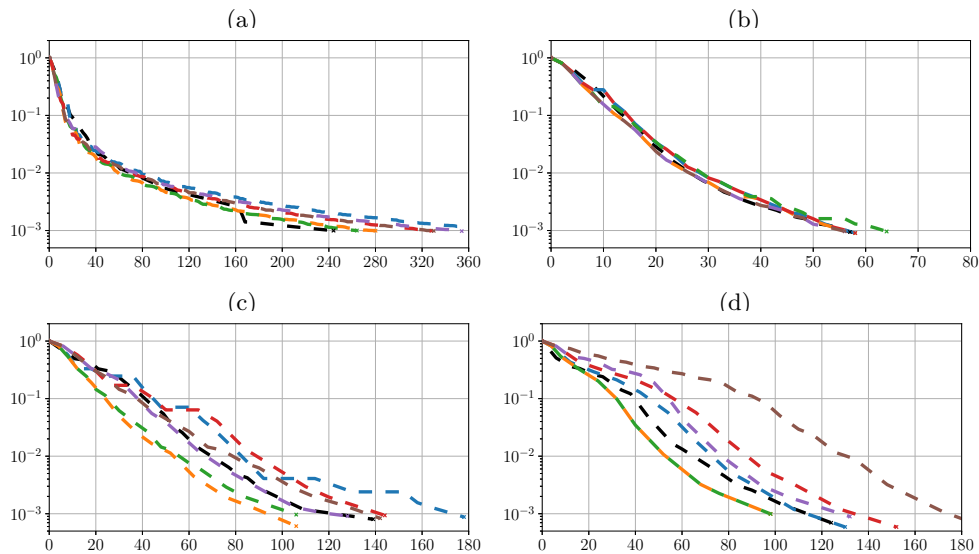


Fig. 5: Data misfit as a function of the computational complexity for the steepest descent (a), the  $l$ -BFGS (b), the full Newton (c) and the Gauss-Newton (d) methods combined with either a line search (—) or a prospective trust region (A (—), B (—), C (—)) or a retrospective trust region (A (—), B (—), C (—)).

671 to anticipate and prevent failures. However the prospective Newton method combined  
 672 with the more cautious parameters sets (B and C) does already not reject any direction  
 673 and there is then no interest in computing the retrospective ratio.

674 Among the prospective methods, it appears that the two more cautious (B and  
 675 C) yield the fastest convergence. Convergence speed decreases when using parameter  
 676 set A with both the full Newton method and the Gauss-Newton method but for two  
 677 different reasons. With parameter set A, the trust-region radius grows quickly and the  
 678 full Newton method is thus allowed to explore large directions, beyond the validity of  
 679 the exact second order expansion (2.2). Such directions produce a high rejection rate  
 680 (32%) and thus a waste of computational effort. At the opposite, the Gauss-Newton  
 681 method never rejects a direction and the explanation for its slower convergence can  
 682 therefore not be the same. During the earliest iterations, far from the global minimum,  
 683 the Gauss-Newton approximation is not valid (because data residuals are not small  
 684 yet) and thus the Gauss-Newton Hessian is quite different from the full Hessian. The  
 685 misfit prediction under the Gauss-Newton approximation is thus cruder than the exact  
 686 second order expansion (2.2) and the ratio  $\rho_p$  is even more likely to be away from  
 687 one. This ratio  $\rho_p$  is given in Fig. 6c. As can be seen, during the first few iterations,  
 688 this ratio is actually very larger than one, which indicates that the misfit prediction  
 689 is indeed not accurate. Nevertheless, the trust region radius is still increased and the  
 690 system is solved accurately while the Hessian and the misfit are not approximated  
 691 accurately. This effect generates over-solving the system at the earliest iteration and  
 692 slows down the Gauss-Newton method, as can be seen by comparing the initial slopes  
 693 between variant A and B/C in Fig. 5d. This effect would be even more dominant  
 694 if the convergence requirements, *i.e.* the forcing sequence  $\eta$ , was smaller. With the  
 695 large value  $\eta = 0.5$  chosen here, convergence of the conjugate gradient algorithm is  
 696 attained relatively fast. Actually variant B and C perform better than variant A only  
 697 because it takes more iterations for the trust region constraint to become inactive.  
 698 Starting with a larger initial radius would result in the same convergence speed than  
 699 variant A. Also, it is interesting to highlight that when using the retrospective radius

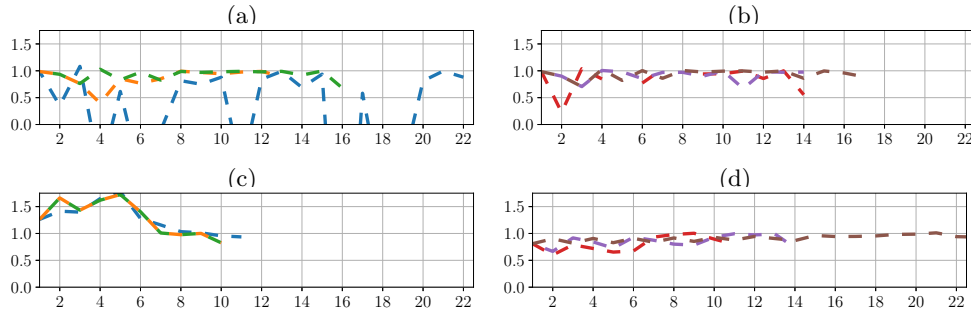


Fig. 6: Prospective ratio  $\rho_p$  (a,c) or retrospective ratio  $\rho_r$  (b,d) during the outer iterations of the full Newton method (a,b) and the Gauss-Newton method (c,d) with different parameter sets using a prospective radius update (a,c) (A (—), B (—), C (—)) or a retrospective radius update (b,d) (A (—), B (—), C (—)).

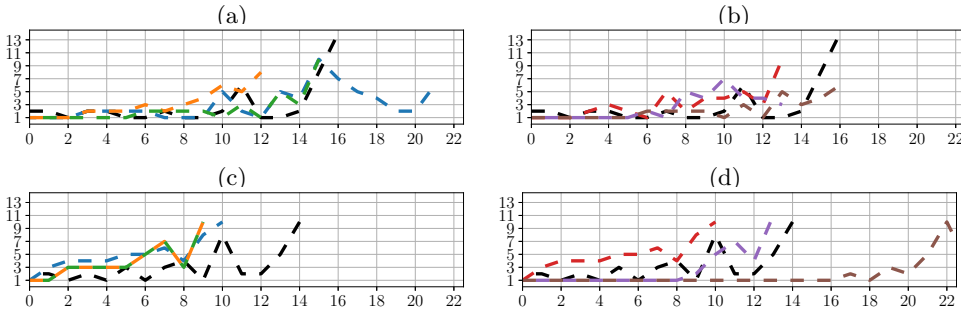


Fig. 7: Inner iterations per outer iteration for the full Newton method (a,b) and the Gauss-Newton method (c,d) with different parameter sets using a prospective radius update (a,c) (A (—), B (—), C (—)) or a retrospective radius update (b,d) (A (—), B (—), C (—)).

700 update with the Gauss-Newton approximation, the situation is reversed because the  
 701 retrospective ratio is then smaller than one. Instead of over-solving, under-solving  
 702 then appears. Therefore we believe that it is better to use trust-region methods with  
 703 the full Newton Hessian, because it constructs the best possible misfit prediction while  
 704 it does not introduce supplementary difficulties.

705 The full Newton method and the Gauss-Newton method are slightly slower when  
 706 combined with a line search method. As far as the full Newton method is concerned,  
 707 directions of very small curvature can produce large update directions, far beyond  
 708 the validity of the expansion (2.2). In such cases the initial length  $\gamma = 1$  is rejected  
 709 and some computational cost must be involved to reduce it to satisfy Wolfe condi-  
 710 tions. This effect has actually been observed twice using the full Newton method.  
 711 Moreover during the first fifth outer iterations, the full Newton method using the line  
 712 search globalization stops because a direction of negative curvature is encountered.  
 713 At the opposite of its trust-region counterpart, the line search variant of the conjugate  
 714 gradient algorithm discard any direction of negative curvature, thus wasting the as-  
 715 sociated computational cost. Of course within the Gauss-Newton approximation this  
 716 second effect can not appear (and the first one was actually not observed). The line  
 717 search globalization therefore seems more suited with the Gauss-Newton approxima-  
 718 tion. Nevertheless it is not much faster. In the context of line search globalization,  
 719 the accuracy of the second order local expansion is expressed through the forcing  
 720 sequence  $\eta$ , which is, as can be seen in Fig. 8, away from zero during the first few iter-  
 721 ations. Consequently, the convergence of the conjugate gradient algorithm is quickly

722 reached and only a few inner iterations are performed per outer iterations as can be  
 723 seen from Fig. 7c. Fig. 7c actually show how hard it is to stop the Gauss-Newton  
 724 inner iterations at the right time: the fastest method is the prospective trust region  
 725 B/C and it performs less inner iterations then the variant A but more than the line  
 726 search method. The difficulty to pick up an appropriate stopping criterion for the  
 727 Gauss-Newton method is another motivation to use the full Newton method instead.  
 728 Using the full Newton method, the line search variant actually suffers from directions  
 729 of small or negative curvature while trust-region methods do not. Based on this case  
 730 study, we would therefore recommend to use the full Newton method combined with  
 731 a trust-region method and a prospective radius update.

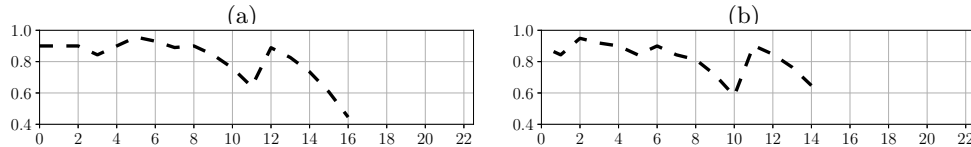


Fig. 8: Forcing sequence  $\eta$  for the full Newton (a) and the Gauss-Newton (b) methods combined with a line search method (—). The forcing sequence for methods combined with a trust-region is constant ( $\eta = 0.5$ ).

732 **3.2. Case study 2.** The configuration of this second case study is replicated  
 733 from [20]. The true velocity distribution is given in Fig. 9a. It presents two T-shaped  
 734 concrete structures ( $v_c = 4$  [km/s]) embedded in a homogeneous background ( $v_b = 0.3$   
 735 [km/s]) with a horizontal layer reflector in the bottom ( $v_r = 0.5$  [km/s]). The depth  
 736 of investigation is limited to 3 [m] while the width is 30 [m]. The aspect ratio and  
 737 the propagation scales are thus very different from the Marmousi model. These two  
 738 concrete foundations, buried at few meters deep, generate high-amplitude reflections  
 739 because of the very high velocity contrast with the background. Moreover, important  
 740 multiple scattering appears between the two structures, as they are relatively close to  
 741 each other. The acquisition system is divided into three segments: one on the surface  
 742 and the two others inside boreholes on both lateral sides. Sources and receivers are  
 743 equally spaced (15 [cm]) along these three segments. Note that the surface sources  
 744 and receivers that would lie inside the two concrete structures are not considered  
 745 in the modelling, leading to an actual number of sources and receivers totaling 227.  
 746 Nine frequencies (100, 125, 150, 175, 200, 225, 250, 275, and 300 [Hz]) are inverted  
 747 simultaneously from an initial model composed of the homogeneous background and  
 748 the bottom reflector only. For this second case study, a logarithmic slowness squared  
 749 parametrization is used  $m := \ln s^2$ . This parametrization has the advantage to be  
 750 unable to produce negative values of the slowness squared. Inverting the slowness  
 751 squared actually drives it into negative values, because of the two concrete structures  
 752 whose slowness squared is really close to zero. Outer iterations are stopped when  
 753 satisfying the convergence criterion  $J(\ln s^2)/J(\ln s_{\text{init}}^2) < 10^{-2}$ . Slowness squared  
 754 fields and pressure fields at the nine frequencies are discretized on a square grid (15  
 755 [cm]) by hierarchical finite elements, respectively of order 1 and of order 2, 3, 4. At  
 756 the light of the first case study, trust-region methods with parameter sets A and C  
 757 will no longer be considered, as both were systematically outperformed by parameter  
 758 set B.

759 **3.2.1. Inner product.** Illumination of the medium is nearly perfect and conse-  
 760 quently, the diagonal part of the Gauss-Newton Hessian that we previously used as a  
 761 weight can reasonably be approximated by a constant  $h_{\text{GN}}$ . However the part related  
 762 to the change of variable is varying spatially  $\delta s^2 = \frac{ds^2}{d \ln s^2} \delta \ln s^2 = s^2 \delta \ln s^2$ . Hence the  
 763 weight for the inner product is chosen as  $w = h_{\text{GN}} s^4$ . As previously, the line search

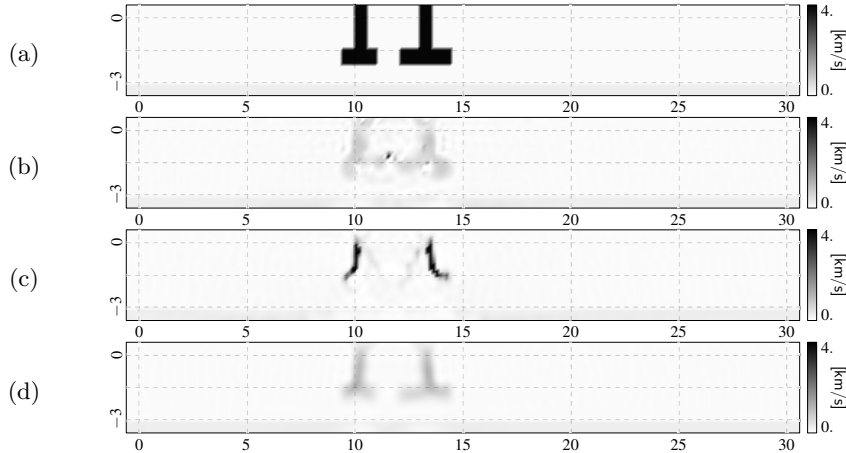


Fig. 9: Near-surface concrete structures velocity model (a) and inversion results using a line search  $l$ -BFGS algorithm with a conventional (b), a weighted only (c) or a weighted and smoothed (d) inner product.

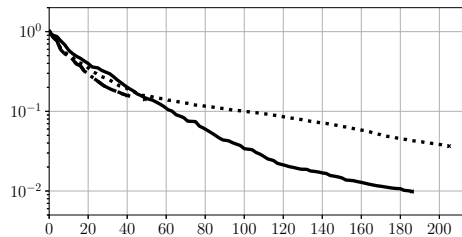


Fig. 10: Data misfit as a function of the computational complexity for the line search  $l$ -BFGS algorithm with a conventional ( $\bullet\bullet$ ), only weighted ( $-\bullet$ ) or weighted and smoothed ( $-$ ) inner product.

764  $l$ -BFGS algorithm has been applied with the four different inner products introduced  
765 in this work. Convergence curves are given in Fig. 10 while inversion results are given  
766 in Fig. 9. For the weighted and smoothed variant, the threshold is set as  $\epsilon = h_{\text{GN}} s_b^4$   
767 while the characteristic length for the smoothing inner product is set to  $l_c = 3$  [m]. It  
768 is important to highlight that this length is greater than the smallest wavelength in  
769 the background medium (1 [m]) while for the first case study, this length was actually  
770 close to the smallest wavelength. The weighted and thresholded variant has been  
771 tested for several values of the threshold, from  $\epsilon = h_{\text{GN}} s_c^4$  to  $\epsilon = h_{\text{GN}} s_b^4$  but none  
772 of them provided inversion results significantly different from the conventional or the  
773 weighted inner products. Only the smoothing inner product is able to reconstruct the  
774 model parameter accurately. This smoothing inner product actually mitigates the  
775 non-linearity of the misfit, because spatial roughness is incorporated progressively in  
776 the model parameter [41]. During the inversion, the model parameter never explores  
777 extremely high velocity values, at the opposite of the other variants. It is thus able  
778 to converge to an accurate solution while more straightforward optimization is not.  
779 Consequently, this inner product is used for the remainder of this study. The perfor-  
780 mance of the three optimization methods is described in the next three subsections.  
781 Convergence curves, inversion results and statistics are given in Fig. 12 and 11 and in  
782 Table 3 respectively.

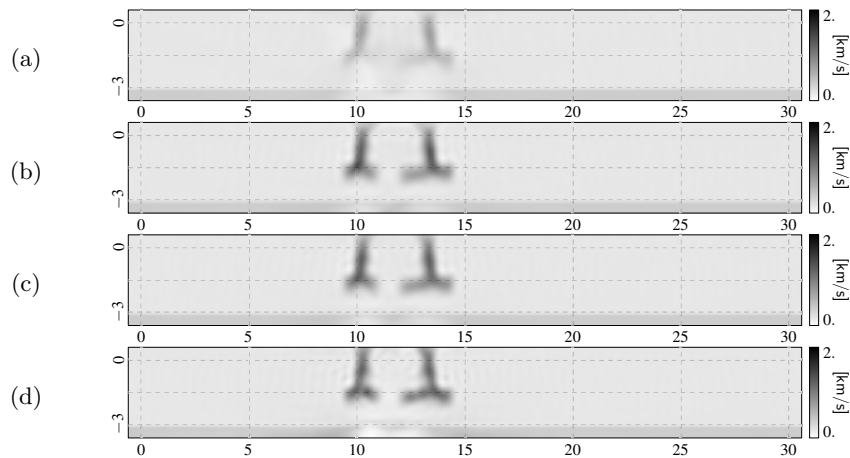


Fig. 11: Inversion results for the steepest descent (a), the  $l$ -BFGS (b), the full Newton (c) and the Gauss-Newton (d) methods combined with trust-region method using a prospective radius update (B). Note the the upper color scale limit is only 2 [km/s].

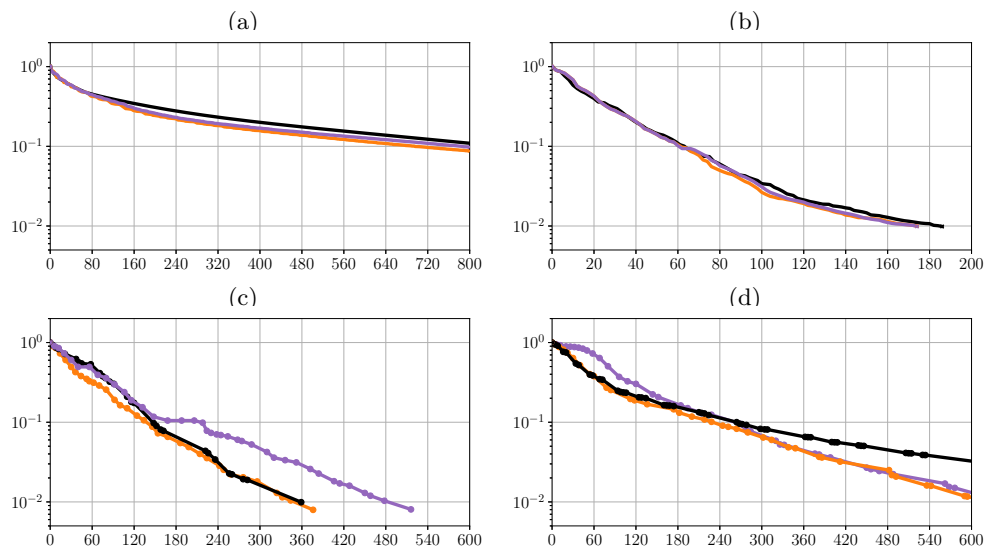


Fig. 12: Data misfit as a function of the computational complexity for the steepest descent (a), the  $l$ -BFGS (b), the full Newton (c) and the Gauss-Newton (d) methods combined with either a line search (—) or a trust-region with a prospective (B —) or a retrospective (B —) radius update. Dots on (Gauss)-Newton curves indicate outer iterations.

783 **3.2.2. Steepest descent.** The steepest descent method is not able to reach  
 784 convergence in a reasonable amount of computations. Progressively decreasing the  
 785 smoothing length  $l_c$  during the inversion would accelerate the convergence [41], but  
 786 it is not needed for more sophisticated methods and thus it is not done here neither.  
 787 As for the first test case, the slope of trust-region methods is slightly steeper than the  
 788 line search method. The prospective radius update rejects less often directions and  
 789 hence converges faster than the retrospective radius update.

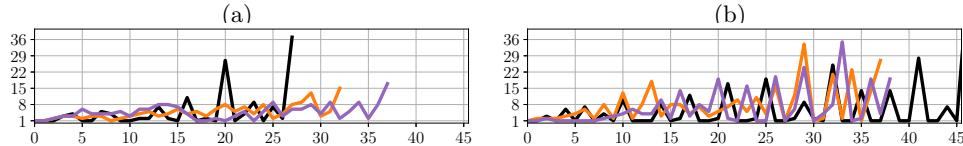
		Wave sol. (tot)	Outer it. (tot)	Inner it. (avg)	Rejected (%)	Constrained (%)	Negative curv. (%)
SD	LS*	803	400	-	1	-	-
	TR-P (B)*	800	400	-	9	100	-
	TR-R (B)*	800	400	-	18	100	-
LB	LS	186	88	-	8	-	-
	TR-P (B)	174	87	-	1	23	-
	TR-R (B)	174	87	-	2	6	-
FN	LS	359	28	5.0	29	-	25
	TR-P (B)	376	33	4.7	0	70	0
	TR-R (B)	516	38	4.8	8	68	0
GN	LS*	923	60	6.6	8	-	-
	TR-P (B)	680	38	7.9	0	42	-
	TR-R (B)	672	39	6.6	0	41	-

Table 3: Statistics related to the implementation of the steepest descent (SD), the  $l$ -BFGS (LB), the full Newton (FN) and the Gauss-newton (GN) methods combined with a line search (LS) or a trust-region (TR) with a prospective (P) or retrospective (R) radius update with parameter set B. Star marker \* indicates methods that have been stopped before convergence.

790 **3.2.3. Limited memory BFGS method.** Similarly to the first test case, the  
791 influence of the globalization method on the convergence speed is small. Trust-region  
792 methods actually spare a part of the line search cost, but it already represents only  
793 a tiny fraction (20 wave solutions) of the overall computational cost (186 wave solu-  
794 tions). Retrospective ratio is again always very close to one and the only difference  
795 between retrospective and prospective radius update is the frequency the size con-  
796 straint is active, although it does not influence the convergence speed.

797 **3.2.4. Newton methods.** For this case study, the full Newton method clearly  
798 outperforms the Gauss-Newton method, independently of the globalization method  
799 used. On the one hand, the convergence speed is much higher and on the other  
800 hand the accuracy of the inversion results is superior. As demonstrated in [20], the  
801 missing negative definite part of the Hessian can prevent the Gauss-Newton method  
802 from reaching convergence. Here, thanks to the inner product preconditioning, every  
803 method is able to find the minimum but the invalidity of the Gauss-Newton approx-  
804 imation impacts the convergence speed and the inversion results. Interestingly, for  
805 the Gauss-Newton method, the retrospective radius update succeeds to compensate  
806 its cost (2 wave solutions per outer iteration). Indeed, during the earliest outer iter-  
807 ations when the Gauss-Newton and the full Hessian are different, we observed that  
808 the retrospective ratio is smaller than one while the prospective ratio is bigger than  
809 one. Consequently the retrospective method performs less inner iterations per outer  
810 iterations than the prospective method (Fig. 13b), and thus avoids early over-solving.  
811 In the end both methods still converge at the same speed, but the retrospective  
812 method has spent less time in the computation of linear system solutions (680 versus  
813  $672 - 2 \times 39 = 594$  wave solutions). At the opposite, for the full Newton method,  
814 the retrospective method spent even more time in the computation of linear system  
815 solutions than the prospective method. The prospective method is actually already  
816 efficient because the prospective misfit prediction is accurate. The line search glob-  
817 alization also provides fast convergence in this case, despite the fact that direc-  
818 tions of negative curvature are often encountered (14 wasted wave solutions) and that the

819 unit step length is often rejected. However the flow of the method is very different  
 820 from trust-region methods. Indeed line search methods have a tendency to compute  
 821 a single very accurate system solution, followed by several very inaccurate system  
 822 solutions as can be seen from Fig. 13a and from the dots spacing in Fig. 12c while  
 823 trust-region methods perform a nearly steadily increasing number of inner iterations  
 824 per outer iteration. Whether a flow is better than the other has not been emphasized  
 825 by our case studies. In the case of noisy data, we however believe it could have an  
 influence.



826 Fig. 13: Inner iterations per outer iteration for the full Newton method (a) and the  
 Gauss-Newton method (b) combined with either a line search (—) or a trust-region  
 with a prospective (B (—)) or a retrospective (B (—)) radius update.

827 **4. Conclusion.** In this work, we investigated the use of trust-region methods  
 828 in the context of full waveform inversion in the frequency domain. At the heart of  
 829 any trust-region method is the trust-region constraint, which is expressed in terms of  
 830 the inner product chosen for the model parameter space. Consequently we begun our  
 831 analysis by investigating different inner product choices that could be implemented.  
 832 We showed that changing the inner product does not only modify how lengths are  
 833 measured but also acts as a preconditioner on both the gradient and the Hessian  
 834 operator. Based on two numerical case studies, we showed that moving from a con-  
 835 ventional inner product to a smoothed and/or weighted inner product can accelerate  
 836 the convergence and mitigate the non-linearity of the misfit, for any optimization  
 837 method independently of the globalization method (line search or trust region).

838 In parallel with this inner product choice, we also introduced line search and  
 839 trust-region variants of the steepest descent, the  $l$ -BFGS and the (Gauss-)Newton  
 840 methods. The number of wave propagation problems to be solved for each method  
 841 was derived in order to compare them fairly. For each optimization method, the line  
 842 search and the trust-region globalizations were then compared based on two different  
 843 case studies. Thanks to the inner product preconditioning, every combination actu-  
 844 ally already yields very satisfying results. Nevertheless, we showed that trust-region  
 845 methods outperform line search methods in numerous situations. In particular, we  
 846 observed that the steepest descent converges slightly faster, because the trust-region  
 847 methods always tried to increase the step length. As far as the  $l$ -BFGS method is  
 848 concerned, very few differences were noted, but interestingly, constraining the size of  
 849 the update direction did not decrease the convergence speed. The more dramatic dif-  
 850 ferences appeared when using the full Newton method. Trust-region methods actually  
 851 overcome the difficulties that appeared when using a line search method with the full  
 852 Newton method. The Gauss-Newton approximation is not required with trust-region  
 853 methods and actually degrades their performances, because this approximation also  
 854 degrades the misfit prediction.

855 We believe that more sophisticated optimization methods, for example combining  
 856  $l$ -BFGS and Newton methods, could increase even more the convergence speed. Future  
 857 works should also investigate the behaviour of inner product preconditioned trust-  
 858 region methods in the presence of noise, possibly with new inner products that involve  
 859 prior information on the model parameter space. We believe that the size constraint  
 860 could act as a regularization method *per se*. Based on our study and these potential  
 861 extensions, trust-region methods and inner product preconditioning seem to be two  
 862 very useful tools for full waveform inversion.



863 **5. Acknowledgements.** The authors would like to thank Anthony Royer for  
 864 his help on the finite element solver used in this work [30].

865

## REFERENCES

- 866 [1] X. ADRIAENS, F. HENROTTE, AND C. GEUZAINÉ, *Adjoint state method for time-harmonic scat-*  
 867 *tering problems with boundary perturbations*, Journal of Computational Physics, 428 (2021),  
 868 p. 109981.
- 869 [2] A. Y. ANAGAW AND M. D. SACCHI, *Model parametrization strategies for Newton-based acoustic*  
 870 *full waveform inversion*, Journal of Applied Geophysics, 157 (2018), pp. 23–36.
- 871 [3] R. BROSSIER, S. OPERTO, AND J. VIRIEUX, *Seismic imaging of complex onshore struc-*  
 872 *tures by 2D elastic frequency-domain full-waveform inversion*, GEOPHYSICS, 74 (2009),  
 873 pp. WCC105–WCC118.
- 874 [4] R. BROSSIER, S. OPERTO, AND J. VIRIEUX, *Which data residual norm for robust elastic*  
 875 *frequency-domain full waveform inversion?*, GEOPHYSICS, 75 (2010), pp. R37–R46.
- 876 [5] M. D. S. R. CARNEIRO, B. PEREIRA-DIAS, D. M. SOARES FILHO, AND L. LANDAU, *On the Scaling*  
 877 *of the Update Direction for Multi-parameter Full Waveform Inversion: Applications to 2D*  
 878 *Acoustic and Elastic Cases*, Pure and Applied Geophysics, 175 (2018), pp. 217–241.
- 879 [6] E. CAUSSE, R. MITTET, AND B. URSIN, *Preconditioning of full-waveform inversion in viscoa-*  
 880 *coustic media*, GEOPHYSICS, 64 (1999), pp. 130–145.
- 881 [7] B. CONSOLVO, M. ZUBERI, R. PRATT, AND P. CARY, *FWI with Scaled-Sobolev Preconditioning*  
 882 *Applied to Short-offset Vibroseis Field Data*, in 79th EAGE Conference and Exhibition  
 883 2017, no. July, 6 2017, pp. 10–15.
- 884 [8] D. DATTA AND M. K. SEN, *Estimating a starting model for full-waveform inversion using a*  
 885 *global optimization method*, GEOPHYSICS, 81 (2016), pp. R211–R223.
- 886 [9] S. C. EISENSTAT AND H. F. WALKER, *Choosing the forcing terms in an inexact Newton method*,  
 887 SIAM Journal of Scientific Computing, 17 (1996), pp. 16–32.
- 888 [10] J. FAN, J. PAN, AND H. SONG, *A Retrospective Trust Region Algorithm with Trust Region*  
 889 *Converging to Zero*, Journal of Computational Mathematics, 34 (2016), pp. 421–436.
- 890 [11] J.-Y. FAN AND Y.-X. YUAN, *A new trust region algorithm with trust region radius converging*  
 891 *to zero*, Proceedings of the 5th International Conference on Optimization: Techniques and  
 892 Applications, (2001).
- 893 [12] Y. FAVENNEC, F. DUBOT, D. LE HARDY, B. ROUSSEAU, AND D. R. ROUSSE, *Space-Dependent*  
 894 *Sobolev Gradients as a Regularization for Inverse Radiative Transfer Problems*, Mathe-  
 895 *matical Problems in Engineering*, 2016 (2016), pp. 1–18.
- 896 [13] A. FICHTNER AND J. TRAMPERT, *Hessian kernels of seismic data functionals based upon adjoint*  
 897 *techniques*, Geophysical Journal International, 185 (2011), pp. 775–798.
- 898 [14] A. GUITTON, G. AYENI, AND E. DÍAZ, *Constrained full-waveform inversion by model reparam-*  
 899 *eterization*, GEOPHYSICS, 77 (2012), pp. R117–R127.
- 900 [15] M. HINZE, R. PINNAU, M. ULBRICH, AND S. ULBRICH, *Optimization with PDE Constraints*,  
 901 vol. 23 of Mathematical Modelling: Theory and Applications, Springer Netherlands, Dor-  
 902 drecht, 2009.
- 903 [16] W. HU, A. ABUBAKAR, T. M. HABASHY, AND J. LIU, *Preconditioned non-linear conjugate gra-*  
 904 *dient method for frequency domain full-waveform seismic inversion*, Geophysical Prospect-  
 905 ing, 59 (2011), pp. 477–491.
- 906 [17] P. KAZEMI AND R. RENKA, *Minimization of the Ginzburg–Landau energy functional by a*  
 907 *Sobolev gradient trust-region method*, Applied Mathematics and Computation, 219 (2013),  
 908 pp. 5936–5942.
- 909 [18] P. LIN, S. PENG, Y. LU, AND W. DU, *The trust region method for time-domain full waveform*  
 910 *inversion*, in Proceedings of the 7th International Conference on Environment and Engi-  
 911 *neering Geophysics & Summit Forum of Chinese Academy of Engineering on Engineering*  
 912 *Science and Technology*, Paris, France, 2016, Atlantis Press, pp. 220–223.
- 913 [19] L. MÉTIVIER, F. BRETAEU, R. BROSSIER, S. OPERTO, AND J. VIRIEUX, *Full waveform*  
 914 *inversion and the truncated Newton method: quantitative imaging of complex subsurface*  
 915 *structures*, Geophysical Prospecting, 62 (2014), pp. 1353–1375.
- 916 [20] L. MÉTIVIER, R. BROSSIER, J. VIRIEUX, AND S. OPERTO, *Full Waveform Inversion and the*  
 917 *Truncated Newton Method*, SIAM Journal on Scientific Computing, 35 (2013), pp. B401–  
 918 B437.
- 919 [21] W. MULDER AND R.-E. PLESSIX, *Exploring some issues in acoustic full waveform inversion*,  
 920 Geophysical Prospecting, 56 (2008), pp. 827–841.
- 921 [22] J. NEUBERGER, *Sobolev Gradients and Differential Equations*, vol. 1670 of Lecture Notes in  
 922 Mathematics, Springer Berlin Heidelberg, Berlin, Heidelberg, 2010.
- 923 [23] J. NOCEDAL, S. J. WRIGHT, AND S. M. ROBINSON, *Numerical Optimization*, Springer Series in  
 924 Operations Research and Financial Engineering, Springer New York, 2006.
- 925 [24] S. OPERTO, Y. GHOLAMI, V. PRIEUX, A. RIBODETTI, R. BROSSIER, L. METIVIER, AND  
 926 J. VIRIEUX, *A guided tour of multiparameter full-waveform inversion with multicompo-*

- 927                    *ment data: From theory to practice*, The Leading Edge, 32 (2013), pp. 1040–1054.
- 928 [25] W. PAN, K. A. INNANEN, AND W. LIAO, *Accelerating Hessian-free Gauss-Newton full-waveform*  
929 *inversion via l-BFGS preconditioned conjugate-gradient algorithm*, GEOPHYSICS, 82  
930 (2017), pp. R49–R64.
- 931 [26] W. PAN, K. A. INNANEN, G. F. MARGRAVE, M. C. FEHLER, X. FANG, AND J. LI, *Estimation*  
932 *of elastic constants for HTI media using Gauss-Newton and full-Newton multiparameter*  
933 *full-waveform inversion*, GEOPHYSICS, 81 (2016), pp. R275–R291.
- 934 [27] B. PARK, W. HA, AND C. SHIN, *A comparison of the preconditioning effects of different param-*  
935 *eterization methods for monoparameter full waveform inversions in the Laplace domain*,  
936 *Journal of Applied Geophysics*, 172 (2020), p. 103883.
- 937 [28] R.-E. PLESSIX, *A review of the adjoint-state method for computing the gradient of a functional*  
938 *with geophysical applications*, *Geophysical Journal International*, 167 (2006), pp. 495–503.
- 939 [29] R. G. PRATT AND C. SHIN, *Gauss-Newton and full Newton methods in frequency-space seismic*  
940 *waveform inversion*, *Geophysical Journal International*, 133 (1998), pp. 341–362.
- 941 [30] A. ROYER, E. BÉCHET, AND C. GEUZAINÉ, *GMSH-FEM : an efficient finite element library*  
942 *based on GMSH*, in 14th World Congress on Computational Mechanics (WCCM), no. July,  
943 2020, pp. 19–24.
- 944 [31] S. H. SCHOT, *Eighty years of Sommerfeld’s radiation condition*, *Historia Mathematica*, 19  
945 (1992), pp. 385–401.
- 946 [32] L. SIRGUE AND R. G. PRATT, *Efficient waveform inversion and imaging: A strategy for selecting*  
947 *temporal frequencies*, GEOPHYSICS, 69 (2004), pp. 231–248.
- 948 [33] T. STEihaug, *The Conjugate Gradient Method and Trust Regions in Large Scale Optimization*,  
949 *SIAM Journal on Numerical Analysis*, 20 (1983), pp. 626–637.
- 950 [34] R. VERSTEEG, *The Marmousi experience: Velocity model determination on a synthetic complex*  
951 *data set*, The Leading Edge, 13 (1994), pp. 927–936.
- 952 [35] J. VIRIEUX AND S. OPERTO, *An overview of full-waveform inversion in exploration geophysics*,  
953 *GEOPHYSICS*, 74 (2009), pp. WCC1–WCC26.
- 954 [36] P. WELLINGTON\*, R. BROSSIER, S. GARAMBOIS, AND J. VIRIEUX, *Laplacian based Precondi-*  
955 *tioning of FWI: Using prior information from seismic reflection data.*, in SEG Technical  
956 *Program Expanded Abstracts 2015*, Society of Exploration Geophysicists, 8 2015, pp. 1436–  
957 1440.
- 958 [37] X. YAN, Q. HE, AND Y. WANG, *Truncated trust region method for nonlinear inverse prob-*  
959 *lems and application in full-waveform inversion*, *Journal of Computational and Applied*  
960 *Mathematics*, 404 (2022), p. 113896.
- 961 [38] P. YANG, R. BROSSIER, L. MÉTIVIER, J. VIRIEUX, AND W. ZHOU, *A Time-Domain Precon-*  
962 *ditioned Truncated Newton Approach to Visco-acoustic Multiparameter Full Waveform*  
963 *Inversion*, *SIAM Journal on Scientific Computing*, 40 (2018), pp. B1101–B1130.
- 964 [39] H. ZHANG, X. LI, H. SONG, AND S. LIU, *An adaptive subspace trust-region method for frequency-*  
965 *domain seismic full waveform inversion*, *Computers & Geosciences*, 78 (2015), pp. 1–14.
- 966 [40] W. ZHANG AND Y. LI, *Elastic wave full-waveform inversion in the time domain by the trust*  
967 *region method*, *Journal of Applied Geophysics*, 197 (2022), p. 104540.
- 968 [41] M. A. ZUBERI AND R. G. PRATT, *Mitigating nonlinearity in full waveform inversion using*  
969 *scaled-Sobolev pre-conditioning*, *Geophysical Journal International*, 213 (2018), pp. 706–  
970 725.