

# List of Supporting Documents

## A. Derivation of the Alternating Direction Method of Multipliers and IMPULSE

## B. Vectorized SAR oracle

- Supporting Information Figure S1
- Supporting Information Figure S2

## C. FAI Update

- Supporting Information Figure S3

## D. VOP compression

- Supporting Information Figure S4
- Supporting Information Figure S5

## E. Tuning the IMPULSE algorithm

- Supporting Information Figure S6
- Supporting Information Video S1 (rho\_small.avi)
- Supporting Information Video S2 (rho\_best.avi)
- Supporting Information Video S3 (rho\_large.avi)

## F. Flip Angle Maps

- Supporting Information Figure S7

## G. High-Level Overview

- Supporting Information Figure S8

## A. Derivation of the Alternating Direction Method of Multipliers and IMPULSE

We want to solve a problem of the form

Problem 1:

$$\begin{aligned} & \text{MINIMIZE } g(y) + h(z) \\ & \text{SUBJECT TO } y = z \end{aligned}$$

where  $g(y)$  and  $h(z)$  are possibly nonconvex functions.

This is equivalent to solving

Problem 2:

$$\text{MAXIMIZE } v(w) = \min_{y,z} (g(y) + h(z) + w^T(y - z))$$

The reasoning for this equivalency is that for a given  $w$ ,  $v(w)$  will give a lower bound on the optimal value for Problem 1. The maximum of this lower bound over all values of  $w$  will give the solution of Problem 1. The constraint  $y = z$  is implicitly satisfied by solving Problem 2 since if  $y \neq z$  then it will always be possible to find some  $w$  such that  $v(w)$  is infinity. Problem 2 is called the dual problem of Problem 1 (which is called the primal problem) and solving this dual problem is a standard approach for equality constrained optimization.  $y$  and  $z$  are called primal variables and  $w$  is called the dual variable or Lagrange multiplier (hence the name method of multipliers). To solve Problem 2, a gradient ascent approach can be used with the basic idea of computing the gradient (or technically subgradient if nondifferentiable) of the dual cost function with respect to the dual variable  $w$  and iteratively updating  $w$  with a step in the positive gradient direction. The gradient (at iteration  $k$ ) can be computed in a two step manner:

1.  $(y^{(k+1)}, z^{(k+1)}) = \operatorname{argmin}_{y,z} (g(y) + h(z) + w^{(k)T}(y - z))$
2.  $\nabla v(w^{(k)}) = y^{(k+1)} - z^{(k+1)}$

Then the update of  $w$  is

$$w^{(k+1)} = w^{(k)} + \alpha(y^{(k+1)} - z^{(k+1)})$$

where  $\alpha$  is the step size.

Since the expression  $g(y) + h(z) + w^{(k)T}(y - z)$  is completely separable in  $y$  and  $z$ ,  $(y^{(k+1)}, z^{(k+1)}) = \operatorname{argmin}_{y,z} (g(y) + h(z) + w^{(k)T}(y - z))$  can be easily decomposed into  $y^{(k+1)} = \operatorname{argmin}_y (g(y) + w^{(k)T}y)$  and  $z^{(k+1)} = \operatorname{argmin}_z (h(z) + w^{(k)T}z)$ . Unfortunately the algorithm as stated has poor convergence properties. This can be improved by adding a penalty or regularization term to the objective of problem 1:

$$\begin{aligned} & \text{MINIMIZE } g(y) + h(z) + \frac{\rho}{2} \|y - z\|_2^2 \\ & \text{SUBJECT TO } y = z \end{aligned}$$

This is equivalent to the previous problem since when the constraint is met, the penalty is 0. The algorithm is now:

1.  $(y^{(k+1)}, z^{(k+1)}) = \operatorname{argmin}_{y,z} (g(y) + h(z) + w^{(k)T}(y - z) + \frac{\rho}{2} \|y - z\|_2^2)$

$$2. \quad \nabla v(w^{(k)}) = y^{(k+1)} - z^{(k+1)}$$

The optimal step size can be shown to be  $\rho$  so the update of  $w$  is

$$w^{(k+1)} = w^{(k)} + \rho(y^{(k+1)} - z^{(k+1)})$$

Unfortunately, because of the penalty term, the first step is no longer separable with  $y$  and  $z$ . Therefore, the joint minimization in step 1 cannot be accomplished exactly in a sequential way. Nonetheless, by performing the updates sequentially (or in alternating directions) and then iterating convergence can be shown to occur (even though at each iteration the exact gradient is not found so the direction of ascent is not the steepest). The  $\rho$  parameter can be tuned to optimize the convergence rate. So the algorithm is

$$\begin{aligned} y^{(k+1)} &= \operatorname{argmin}_y \left( g(y) + h(z^{(k)}) + w^{(k)T} (y - z^{(k)}) + \frac{\rho}{2} \|y - z^{(k)}\|_2^2 \right) \\ &= \operatorname{argmin}_y \left( g(y) + w^{(k)T} y + \frac{\rho}{2} \|y\|_2^2 - \rho y^T z^{(k)} \right) \\ &= \operatorname{argmin}_y \left( g(y) + \frac{\rho}{2} \|y\|_2^2 + \rho y^T \left( \frac{1}{\rho} w^{(k)} - z^{(k)} \right) \right) \\ &= \operatorname{argmin}_y \left( g(y) + \frac{\rho}{2} \left\| y + \left( \frac{1}{\rho} w^{(k)} - z^{(k)} \right) \right\|_2^2 \right) \end{aligned}$$

$$\begin{aligned} z^{(k+1)} &= \operatorname{argmin}_z \left( g(y^{(k+1)}) + h(z) + w^{(k)T} (y^{(k+1)} - z) + \frac{\rho}{2} \|y^{(k+1)} - z\|_2^2 \right) \\ &= \operatorname{argmin}_z \left( h(z) + u^{(k)T} z + \frac{\rho}{2} \|z\|_2^2 - \rho z^T y^{(k+1)} \right) \\ &= \operatorname{argmin}_z \left( h(z) + \frac{\rho}{2} \|z\|_2^2 + \rho z^T \left( \frac{1}{\rho} w^{(k)} - y^{(k+1)} \right) \right) \\ &= \operatorname{argmin}_z \left( h(z) + \frac{\rho}{2} \left\| z + \left( \frac{1}{\rho} w^{(k)} - y^{(k+1)} \right) \right\|_2^2 \right) \end{aligned}$$

$$w^{(k+1)} = w^{(k)} + \rho(y^{(k+1)} - z^{(k+1)})$$

For simplicity define  $u = \frac{1}{\rho} w$  so the complete ADMM algorithm for each iteration is:

1.  $y^{(k+1)} = \operatorname{argmin}_y \left( g(y) + \frac{\rho}{2} \|y + u^{(k)} - z^{(k)}\|_2^2 \right)$
2.  $z^{(k+1)} = \operatorname{argmin}_z \left( h(z) + \frac{\rho}{2} \|z + u^{(k)} - y^{(k+1)}\|_2^2 \right)$
3.  $u^{(k+1)} = u^{(k)} + y^{(k+1)} - z^{(k+1)}$

This ADMM algorithm can be applied to the minSAR problem by first expressing it as an unconstrained minimization:

$$\text{MINIMIZE: COST}(\underline{\mathbf{B}}) + \sum_{s=1}^{N_S} I_{(s)}(\underline{\mathbf{b}}_{(s)})$$

Where  $I_{(s)}$  is the indicator function for the set  $\mathcal{F}_{(s)}$ :  $I_{(s)}(\underline{\mathbf{b}}_{(s)}) = \begin{cases} 0, & \underline{\mathbf{b}}_{(s)} \in \mathcal{F}_{(s)} \\ \infty, & \underline{\mathbf{b}}_{(s)} \notin \mathcal{F}_{(s)} \end{cases}$ .

(This is a standard notational trick to express a constrained minimization as an unconstrained minimization: if  $\underline{\mathbf{b}}_{(s)} \notin \mathcal{F}_{(s)}$ , the value will be  $\infty$  which is equivalent to a constraint being violated.) Now applying ADMM to minSAR problem is straightforward by defining  $g(\underline{\mathbf{Y}}) = \text{COST}(\underline{\mathbf{Y}})$  and  $h(\underline{\mathbf{Z}}) = \sum_{s=1}^{N_S} I_{(s)}(\underline{\mathbf{z}}_{(s)})$ . Thus the algorithm is:

GIVEN:  $k = 0; \underline{\mathbf{Y}}^{(0)}, \underline{\mathbf{U}}^{(0)}, \underline{\mathbf{Z}}^{(0)}$

REPEAT:

$$\underline{\mathbf{Y}}^{(k+1)} = \underset{\underline{\mathbf{Y}}}{\text{argmin}} (\text{COST}(\underline{\mathbf{Y}}) + \frac{\rho}{2} \|\underline{\mathbf{Y}} - \underline{\mathbf{Z}}^{(k)} + \underline{\mathbf{U}}^{(k)}\|_2^2)$$

$$\underline{\mathbf{Z}}^{(k+1)} = \underset{\underline{\mathbf{Z}}}{\text{argmin}} \left( \sum_{s=1}^{N_S} I_{(s)}(\underline{\mathbf{z}}_{(s)}) + \frac{\rho}{2} \|\underline{\mathbf{Y}}^{(k+1)} - \underline{\mathbf{Z}} + \underline{\mathbf{U}}^{(k)}\|_2^2 \right)$$

$$\underline{\mathbf{U}}^{(k+1)} = \underline{\mathbf{U}}^{(k)} + \underline{\mathbf{Y}}^{(k+1)} - \underline{\mathbf{Z}}^{(k+1)}$$

UNTIL:  $\underline{\mathbf{Z}}^{(k+1)} = \underline{\mathbf{Y}}^{(k+1)}$

The z-update can be further decomposed into  $N_S$  independent (and parallelizable) operations

$$\underline{\mathbf{z}}_{(s)}^{(k+1)} = \underset{\underline{\mathbf{z}}_{(s)}}{\text{argmin}} \left( I_{(s)}(\underline{\mathbf{z}}_{(s)}) + \frac{\rho}{2} \|\underline{\mathbf{y}}_{(s)}^{(k+1)} - \underline{\mathbf{z}}_{(s)} + \underline{\mathbf{u}}_{(s)}^{(k+1)}\|_2^2 \right) \text{ for } s = 1, \dots, N_S$$

This operation minimizes the deviation of the variable  $\underline{\mathbf{z}}$  from  $\underline{\mathbf{y}}_{(s)}^{(k+1)} + \underline{\mathbf{u}}_{(s)}^{(k+1)}$  while requiring that  $\underline{\mathbf{z}}$  be a feasible pulse ( $\underline{\mathbf{z}} \in \mathcal{F}_{(s)}$ ). In fact, the  $\rho$  parameter becomes irrelevant and this step is equivalent to the Euclidean projection of the point

$\underline{\mathbf{y}}_{(s)}^{(k+1)} + \underline{\mathbf{u}}_{(s)}^{(k+1)}$  onto the set  $\mathcal{F}_{(s)}$ :

$$\underline{\mathbf{z}}_{(s)}^{(k+1)} = \text{proj}_{\mathcal{F}_{(s)}}(\underline{\mathbf{y}}_{(s)}^{(k+1)} + \underline{\mathbf{u}}_{(s)}^{(k+1)}) \text{ for } s = 1, \dots, N_S$$

So the complete ADMM algorithm for the minSAR problem is:

INITIALIZE:  $k = 0; \underline{\mathbf{U}}^{(0)}, \underline{\mathbf{Z}}^{(0)}$

REPEAT:

$$\underline{\mathbf{Y}}^{(k+1)} = \underset{\underline{\mathbf{Y}}}{\text{argmin}} (\text{COST}(\underline{\mathbf{Y}}) + \frac{\rho}{2} \|\underline{\mathbf{Y}} - \underline{\mathbf{Z}}^{(k)} + \underline{\mathbf{U}}^{(k)}\|_2^2)$$

(PARALLEL) FOR  $s = 1, \dots, N_S$

$$\underline{\mathbf{z}}_{(s)}^{(k+1)} = \text{proj}_{\mathcal{F}_{(s)}}(\underline{\mathbf{y}}_{(s)}^{(k+1)} + \underline{\mathbf{u}}_{(s)}^{(k+1)})$$

END FOR

$$\underline{\mathbf{U}}^{(k+1)} = \underline{\mathbf{U}}^{(k)} + \underline{\mathbf{Y}}^{(k+1)} - \underline{\mathbf{Z}}^{(k+1)}$$

UNTIL:  $\underline{\mathbf{Z}}^{(k+1)} = \underline{\mathbf{Y}}^{(k+1)}$

## B. Vectorized SAR oracle

We wish to construct an ‘‘oracle’’ that is queried with a sequence of channel weightings  $\underline{\mathbf{B}} \in \mathbb{C}^{N_c \times N_K N_s}$  and returns the value of  $\text{COST}(\underline{\mathbf{B}}) \in \mathbb{R}$ , where  $\text{COST}$  has the form  $\text{COST}(\underline{\mathbf{B}}) = \max_{r=1, \dots, N_R} (\sum_{i=1}^{N_s N_K} \underline{\mathbf{B}}_i^H \mathbf{R}_r \underline{\mathbf{B}}_i)$ , and its gradient (or more correctly its subgradient since  $\text{COST}$  is nondifferentiable)  $\underline{\mathbf{G}}(\underline{\mathbf{B}}) \in \mathbb{C}^{N_c \times N_K N_s}$  where  $\underline{\mathbf{G}}$  satisfies  $\text{COST}(\underline{\mathbf{B}} + \Delta \underline{\mathbf{B}}) \geq \text{COST}(\underline{\mathbf{B}}) + \sum_{i=1}^{N_K N_s} \underline{\mathbf{G}}_i^T \Delta \underline{\mathbf{B}}_i$  for any  $\Delta \underline{\mathbf{B}} \in \mathbb{C}^{N_c \times N_K N_s}$  since  $\text{COST}$  is convex. The oracle makes use of vectorized computation that allows for simple GPU acceleration. While this new technique is especially important for IMPULSE as the oracle must be queried numerous times throughout the course of the optimization, it is fundamentally more efficient than naive methods relying on ‘‘for loops’’, for computing time- and spatially-averaged local SAR or temperature for an entire pulse sequence with arbitrarily many voxels and RF pulses.

The general procedure to compute  $\text{COST}$  and  $\underline{\mathbf{G}}$  for a particular  $\underline{\mathbf{B}}$  (shown in detail in Figure B) is:

1. Generate a vector  $\mathbf{p} \in \mathbb{R}^{N_R}$  where the  $r$ th element of  $\mathbf{p}$  is  $\sum_{i=1}^{N_s N_K} \underline{\mathbf{B}}_i^H \mathbf{R}_r \underline{\mathbf{B}}_i$
2. Find the value and index of the largest element of  $\mathbf{p}$ :  $[\text{COST}, r_{\max}] = \mathbf{max}(\mathbf{p})$
3. Compute a subgradient analytically:  $\underline{\mathbf{G}} = 2\mathbf{R}_{r_{\max}} \underline{\mathbf{B}}$

The most costly step is the first one and we describe a vectorized method to accomplish it efficiently that is orders of magnitude more efficient than using a ‘‘for loop’’. Not only is the vectorized method fundamentally more efficient in terms of total number of floating point operations needed, but because it expresses the computation as a matrix vector multiplication, there is an additional efficiency gain when using mathematical software that has built-in optimizations for matrix multiplication such as MATLAB. An added benefit is that because of MATLAB’s built-in abstraction for using GPUs to multiply matrices (*gpuArray*), GPU acceleration can be achieved with no extra programming effort.

As an overview of the procedure, it is possible to exploit the symmetry of  $\mathbf{R}_r$  to evaluate  $\underline{\mathbf{B}}_i^H \mathbf{R}_r \underline{\mathbf{B}}_i$  in a minimal number of floating point operations as the inner product of two *real* valued vectors  $R_{r,vec}$  and  $B_{i,vec}$ . Then  $\sum_{i=1}^{N_s N_K} \underline{\mathbf{B}}_i^H \mathbf{R}_r \underline{\mathbf{B}}_i$  can be computed efficiently as the inner product of  $R_{r,vec}$  and  $\bar{\mathbf{b}} = \sum_{i=1}^{N_s N_K} B_{i,vec}$ . The entire vector  $\mathbf{p}$  can be constructed as a single matrix multiplication  $\bar{\mathbf{R}}^T \bar{\mathbf{b}}$  where  $\bar{\mathbf{R}}$  is a matrix with the  $r$ th column equal to  $R_{r,vec}$ .

The expressions for  $R_{r,vec}$  and  $B_{i,vec}$  are found by expanding  $\underline{\mathbf{B}}_i^H \mathbf{R}_r \underline{\mathbf{B}}_i$ . For simplicity, we will drop the subscripts and use the notation  $\mathbf{b}^H \mathbf{R} \mathbf{b}$  in the derivation where single subscripts indicate vector elements and double subscripts indicate matrix elements.

$$\begin{aligned} \mathbf{b}^H \mathbf{R} \mathbf{b} &= \sum_{i=1}^{N_c} \sum_{j=1}^{N_c} \mathbf{R}_{ij} \mathbf{b}_i^* \mathbf{b}_j = \sum_{i=1}^{N_c} \mathbf{R}_{ii} |\mathbf{b}_i|^2 + \sum_{i=1}^{N_c} \sum_{j=1}^i (\mathbf{R}_{ij} \mathbf{b}_i^* \mathbf{b}_j + \mathbf{R}_{ij}^* \mathbf{b}_i \mathbf{b}_j^*) = \sum_{i=1}^{N_c} \mathbf{R}_{ii} |\mathbf{b}_i|^2 + 2 \sum_{i=1}^{N_c} \sum_{j=1}^i \text{Re}\{\mathbf{R}_{ij} \mathbf{b}_i^* \mathbf{b}_j\} \\ &= \sum_{i=1}^{N_c} \mathbf{R}_{ii} |\mathbf{b}_i|^2 + 2 \sum_{i=1}^{N_c} \sum_{j=i+1}^{N_c} (\text{Re}\{\mathbf{R}_{ij}\} \text{Re}\{\mathbf{b}_i^* \mathbf{b}_j\} - \text{Im}\{\mathbf{R}_{ij}\} \text{Im}\{\mathbf{b}_i^* \mathbf{b}_j\}) \end{aligned}$$

It can be seen that the above expression can be written as the inner product of two real valued vectors  $\hat{\mathbf{R}}, \hat{\mathbf{b}} \in \mathbb{R}^{N_c^2}$

$$\hat{R} = \begin{bmatrix} \mathbf{R}_{1,1} \\ \vdots \\ \mathbf{R}_{N_c, N_c} \\ 2\text{Re}\{\mathbf{R}_{12}\} \\ -2\text{Im}\{\mathbf{R}_{12}\} \\ \vdots \\ 2\text{Re}\{\mathbf{R}_{1, N_c}\} \\ -2\text{Im}\{\mathbf{R}_{1, N_c}\} \\ \vdots \\ 2\text{Re}\{\mathbf{R}_{N_c-1, N_c}\} \\ -2\text{Im}\{\mathbf{R}_{N_c-1, N_c}\} \end{bmatrix} = \begin{bmatrix} 1 \\ \vdots \\ 1 \\ 2 \\ -2 \\ \vdots \\ 2 \\ -2 \\ \vdots \\ 2 \\ -2 \end{bmatrix} \circ \begin{bmatrix} \mathbf{R}_{1,1} \\ \vdots \\ \mathbf{R}_{N_c, N_c} \\ \text{Re}\{\mathbf{R}_{12}\} \\ \text{Im}\{\mathbf{R}_{12}\} \\ \vdots \\ \text{Re}\{\mathbf{R}_{1, N_c}\} \\ \text{Im}\{\mathbf{R}_{1, N_c}\} \\ \vdots \\ \text{Re}\{\mathbf{R}_{N_c-1, N_c}\} \\ \text{Im}\{\mathbf{R}_{N_c-1, N_c}\} \end{bmatrix} \quad \text{and} \quad \hat{b} = \begin{bmatrix} |\mathbf{b}_1|^2 \\ \vdots \\ |\mathbf{b}_{N_c}|^2 \\ \text{Re}\{\mathbf{b}_1^* \mathbf{b}_2\} \\ \text{Im}\{\mathbf{b}_1^* \mathbf{b}_2\} \\ \vdots \\ \text{Re}\{\mathbf{b}_1^* \mathbf{b}_{N_c}\} \\ \text{Im}\{\mathbf{b}_1^* \mathbf{b}_{N_c}\} \\ \vdots \\ \text{Re}\{\mathbf{b}_{N_c-1}^* \mathbf{b}_{N_c}\} \\ \text{Im}\{\mathbf{b}_{N_c-1}^* \mathbf{b}_{N_c}\} \end{bmatrix} = \begin{bmatrix} \text{Re}\{\mathbf{b}_1\}^2 \\ \vdots \\ \text{Re}\{\mathbf{b}_{N_c}\}^2 \\ \text{Re}\{\mathbf{b}_1\}\text{Re}\{\mathbf{b}_2\} \\ -\text{Im}\{\mathbf{b}_1\}\text{Re}\{\mathbf{b}_2\} \\ \vdots \\ \text{Re}\{\mathbf{b}_1\}\text{Re}\{\mathbf{b}_{N_c}\} \\ -\text{Im}\{\mathbf{b}_1\}\text{Re}\{\mathbf{b}_{N_c}\} \\ \vdots \\ \text{Re}\{\mathbf{b}_{N_c-1}\}\text{Re}\{\mathbf{b}_{N_c}\} \\ -\text{Im}\{\mathbf{b}_{N_c-1}\}\text{Re}\{\mathbf{b}_{N_c}\} \end{bmatrix} +$$

$$\begin{bmatrix} \text{Im}\{\mathbf{b}_1\}^2 \\ \vdots \\ \text{Im}\{\mathbf{b}_{N_c}\}^2 \\ \text{Im}\{\mathbf{b}_1\}\text{Im}\{\mathbf{b}_2\} \\ \text{Re}\{\mathbf{b}_1\}\text{Im}\{\mathbf{b}_2\} \\ \vdots \\ \text{Im}\{\mathbf{b}_1\}\text{Im}\{\mathbf{b}_{N_c}\} \\ \text{Re}\{\mathbf{b}_1\}\text{Im}\{\mathbf{b}_{N_c}\} \\ \vdots \\ \text{Im}\{\mathbf{b}_{N_c-1}\}\text{Im}\{\mathbf{b}_{N_c}\} \\ \text{Re}\{\mathbf{b}_{N_c-1}\}\text{Im}\{\mathbf{b}_{N_c}\} \end{bmatrix} = \begin{bmatrix} 1 \\ \vdots \\ 1 \\ 1 \\ -1 \\ \vdots \\ 1 \\ -1 \\ \vdots \\ 1 \\ -1 \end{bmatrix} \circ \begin{bmatrix} \text{Re}\{\mathbf{b}_1\} \\ \vdots \\ \text{Re}\{\mathbf{b}_{N_c}\} \\ \text{Re}\{\mathbf{b}_1\} \\ \text{Im}\{\mathbf{b}_1\} \\ \vdots \\ \text{Re}\{\mathbf{b}_1\} \\ \text{Im}\{\mathbf{b}_1\} \\ \vdots \\ \text{Re}\{\mathbf{b}_{N_c-1}\} \\ \text{Im}\{\mathbf{b}_{N_c-1}\} \end{bmatrix} \circ \begin{bmatrix} \text{Re}\{\mathbf{b}_1\} \\ \vdots \\ \text{Re}\{\mathbf{b}_{N_c}\} \\ \text{Re}\{\mathbf{b}_2\} \\ -\text{Re}\{\mathbf{b}_2\} \\ \vdots \\ \text{Re}\{\mathbf{b}_{N_c}\} \\ -\text{Re}\{\mathbf{b}_{N_c}\} \\ \vdots \\ \text{Re}\{\mathbf{b}_{N_c}\} \\ -\text{Re}\{\mathbf{b}_{N_c}\} \end{bmatrix} + \begin{bmatrix} \text{Im}\{\mathbf{b}_1\} \\ \vdots \\ \text{Im}\{\mathbf{b}_{N_c}\} \\ \text{Im}\{\mathbf{b}_1\} \\ \text{Re}\{\mathbf{b}_1\} \\ \vdots \\ \text{Im}\{\mathbf{b}_1\} \\ \text{Re}\{\mathbf{b}_1\} \\ \vdots \\ \text{Im}\{\mathbf{b}_{N_c-1}\} \\ \text{Re}\{\mathbf{b}_{N_c-1}\} \end{bmatrix} \circ \begin{bmatrix} \text{Im}\{\mathbf{b}_1\} \\ \vdots \\ \text{Im}\{\mathbf{b}_{N_c}\} \\ \text{Im}\{\mathbf{b}_2\} \\ \text{Im}\{\mathbf{b}_2\} \\ \vdots \\ \text{Im}\{\mathbf{b}_{N_c}\} \\ \text{Im}\{\mathbf{b}_{N_c}\} \\ \vdots \\ \text{Im}\{\mathbf{b}_{N_c}\} \\ \text{Im}\{\mathbf{b}_{N_c}\} \end{bmatrix}$$

In this way it is seen that  $\mathbf{b}^H \mathbf{R} \mathbf{b} = \hat{R}^T \hat{b}$  where  $\hat{R} = \text{sgn}_R \circ \mathbf{R}[\text{idx}_R]$  and  $\hat{b} = \text{sgn}_b \circ b_{[1]} \circ b_{[2]} + b_{[3]} \circ b_{[4]} = \text{sgn}_b \circ \mathbf{b}[\text{idx}_{b,1}] \circ \mathbf{b}[\text{idx}_{b,2}] + \mathbf{b}[\text{idx}_{b,3}] \circ \mathbf{b}[\text{idx}_{b,4}]$ ;  $\circ$  is the elementwise product;  $\text{sgn}_R$  is a  $N_c^2$  element vector consisting of 1,  $-2$  and  $2$ ;  $\text{sgn}_b$  is a  $N_c^2$  element vector consisting of 1, and  $-1$ ;  $\mathbf{R}_{[]} = \mathbf{R}[\text{idx}_R]$ ,  $b_{[1]} = \mathbf{b}[\text{idx}_{b,1}]$ ,  $b_{[2]} = \mathbf{b}[\text{idx}_{[2]}]$ ,  $b_{[3]} = \mathbf{b}[\text{idx}_{[3]}]$ , and  $b_{[4]} = \mathbf{b}[\text{idx}_{[4]}]$  are  $N_c^2$  element real-valued vectors formed by accessing the real or imaginary components of  $\mathbf{R}$  or  $\mathbf{b}$  at the locations specified by  $\text{idx}_R$ ,  $\text{idx}_{b,1}$ ,  $\text{idx}_{b,2}$ ,  $\text{idx}_{b,3}$ , and  $\text{idx}_{b,4}$ .

Returning to the original notation it can be seen that  $\hat{R}$  corresponds to  $R_{r,vec}$  and  $\hat{b}$  corresponds to  $B_{i,vec}$ . Therefore, it can be seen that the matrix  $\bar{R} \in \mathbb{R}^{N_c^2 \times N_R}$  can be constructed with each column computed as above for a particular  $\mathbf{R}_r$ . Each  $B_{i,vec}$  is also computed as above for a particular  $\mathbf{B}_i$ . It is possible to precompute  $\bar{R}$ ,  $\text{sgn}$ ,  $\text{idx}_{[1]}$ ,  $\text{idx}_{[2]}$ ,  $\text{idx}_{[3]}$ , and  $\text{idx}_{[4]}$  offline. Evaluating  $p$  for a particular  $\mathbf{B}$  (step 1 in Supporting Information Figure S1), requires the following subcomputations:

- A. Generate  $B_{[1]}, B_{[2]}, B_{[3]}, B_{[4]} \in \mathbb{R}^{N_c^2 \times N_S N_K}$  from  $\mathbf{B} \in \mathbb{C}^{N_c \times N_S N_K}$ .

This involves constructing each column of  $B_{[1]}, B_{[2]}, B_{[3]}, B_{[4]}$  from each column of  $\mathbf{B}$  by the procedure described above through reindexing at precomputed indices. This procedure requires no floating point computations and negligible time.

- B. Compute  $\bar{B} = \text{SGN} \circ B_{[1]} \circ B_{[2]} + B_{[3]} \circ B_{[4]}$ .

This requires three elementwise matrix products and one matrix addition. The matrix product with SGN (a matrix with  $N_S N_K$  columns each of which is equal to  $\text{sgn}$ ) takes negligible time since it is not a floating point computation – SGN consists of only  $-1$  and  $1$ . Thus this step has a total of 3 floating point matrix computations (2 multiplications, 1 addition). The number of floating point operations (FLOPs) for each is equal to the number of elements in the matrix or  $N_C^2 N_S N_K$  and thus the total number of FLOPs is  $3N_C^2 N_S N_K$ . However, because each of the  $N_C^2 N_S N_K$  FLOPs can be completely parallelized across elements of the matrix with appropriate multithreading, the number of matrix computations (=3) is a better measure of the complexity of this step than the total number of FLOPs.

C. Compute  $\bar{b} = \bar{B}\mathbf{1}$ .

This requires 1 matrix operation (summing across columns of  $\bar{B}$ ) and a total of  $N_C^2 N_S N_K$  FLOPs. With appropriate multithreading, the computational complexity is  $O(\log(N_S N_K))$  (since multithreaded summing of a list of numbers has computation complexity  $O(\log(\text{length of list}))$  and complete parallelization of computation across rows of  $\bar{B}$  is assumed).

D. Compute matrix multiplication  $p = \bar{R}^T \bar{b}$ .

This requires 1 matrix multiplication and a total  $N_R N_C^4$  total FLOPs. The complexity of matrix multiplication with multithreading is difficult to describe but with the most efficient GPU implementations, duration of matrix multiplication can be made to be approximately constant until the size of the matrices exceeds the memory capacity of the GPU

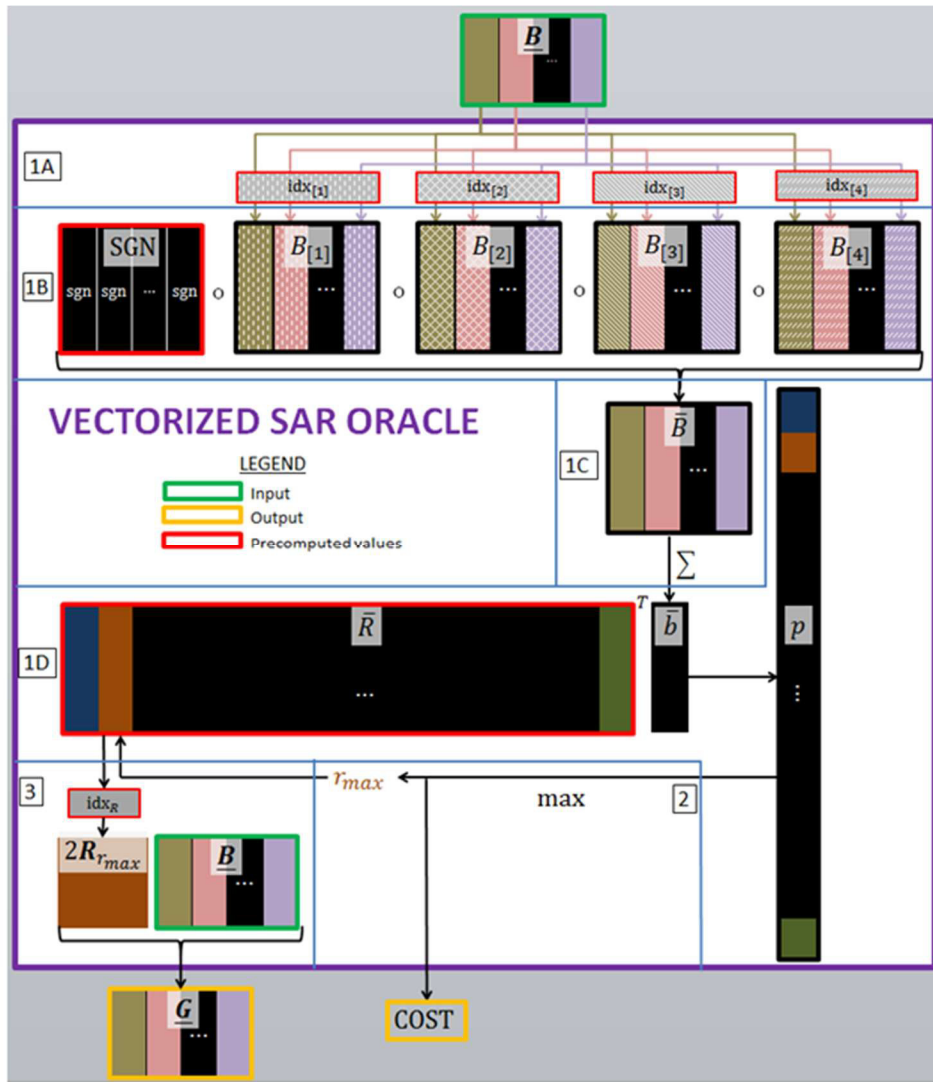
Once the vector  $p$  has been computed, step 2 involves finding the index ( $r_{max}$ ) and value (COST) of the largest element of  $p$ . Step 3 requires selecting the  $r_{max}$ -th column of  $\bar{R}$  and reconstructing the matrix  $\mathbf{R}_{r_{max}}$  through a reverse vectorization operation (which takes negligible time since it only involves reindexing at the predetermined locations,  $\text{idx}_R$ ). Then the subgradient matrix can be found as  $\underline{G} = 2\mathbf{R}_{r_{max}}\underline{B}$ . Both steps 2 and 3 take negligible time compared to step 1. Supporting Information Figure S1 encapsulates this entire procedure graphically.

Based on the analysis above, assuming a very large number of threads (ie with a GPU) the computation time for the oracle is not dependent on the number of floating point operations and takes constant time. The requirement for this is that  $\bar{R}$ , SGN,  $\text{idx}_{[1]}$ ,  $\text{idx}_{[2]}$ ,  $\text{idx}_{[3]}$ ,  $\text{idx}_{[4]}$ , and  $\text{idx}_R$  have been stored onto the GPU prior to the start of optimization. In real time, each input  $\underline{B}$  is transferred to the GPU, steps 1 through 3 are performed on the GPU, then  $\underline{G}$  and COST are transferred back to the CPU. Alternatively, if a GPU is not available, the computations can all be performed on the CPU and significant speedup compared to use of a “for” loop will still be observed. However, memory transfer from CPU to GPU can take significant amount of time and thus it is desirable to perform as much computation on the CPU (at the expense of incomplete parallelization) in order to reduce the amount of data that must be transferred to the GPU in real-time. It is necessary to devise a strategy that manages this tradeoff by deciding when to transfer data and begin performing computations on the GPU (the input  $\underline{B}$  is found on the CPU as an output of the previous ADMM iteration). One option is to transfer  $\underline{B}$  to the GPU and perform Steps 1 through 4 on the GPU; another option is to perform steps 1 through 3 on the CPU and then transfer  $b_{full}$  to the GPU for step 4 (in both cases  $R_{full}$  is assumed to already be stored on the GPU). The advantage of the first option is that steps 2 and 3 are better parallelized on a GPU and will therefore take less time than in the second option. However, the advantage of the second option is that transferring the entire large matrix  $\underline{B}$  is undesirable

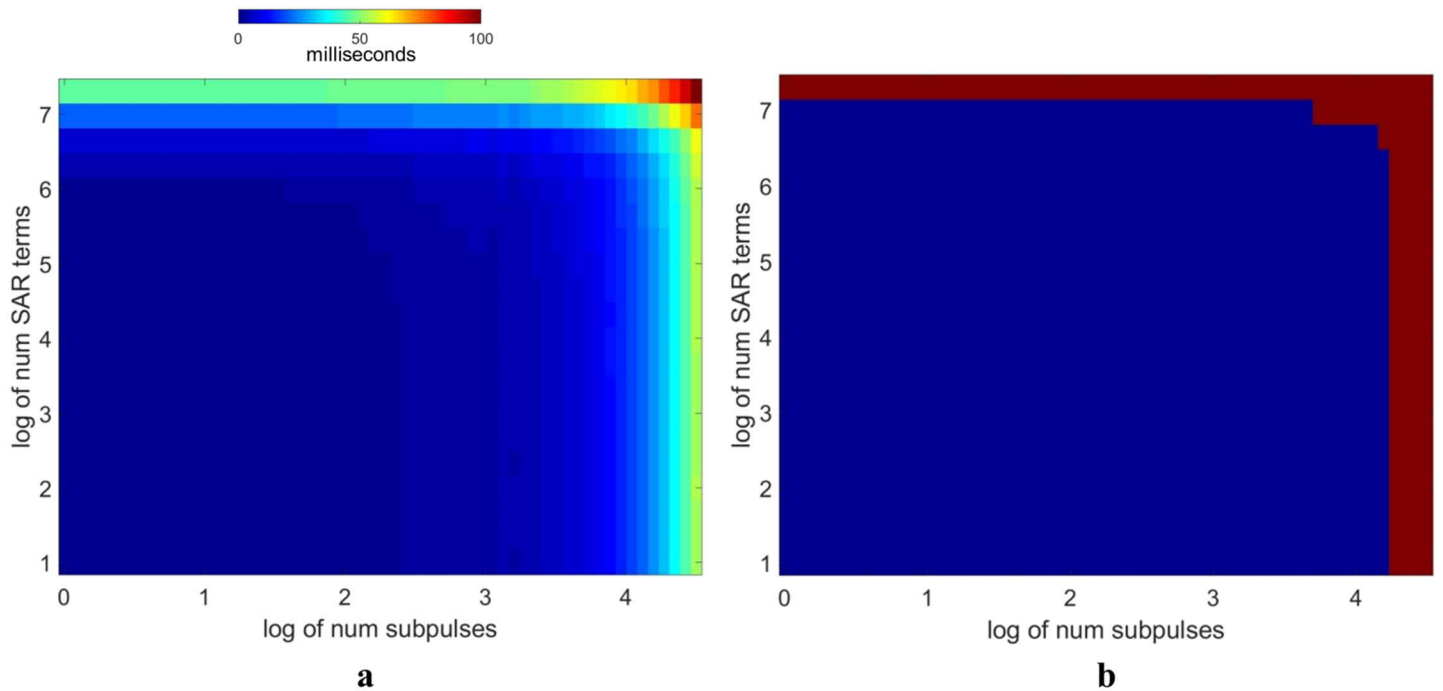
and time-consuming compared to transferring the much smaller vector  $b_{full}$ . Additionally, GPUs have the disadvantage that the reindexing operation in step 1 which takes negligible time on a CPU actually takes a significant amount of time on a GPU. The performance of each option will vary depending on the exact specifications of the processors used for computation.

In summary, by completely vectorizing the computation of SAR and its subgradient it is possible to exploit built-in optimizations from matrix multiplications in most numerical computing packages. Describing the complexity of the vectorized computation exactly is difficult because of the intricacies of the caching and single input multiple data (SIMD) optimizations made by the numerical computing package but empirically favorable scaling is observed with both number of SAR terms and number of pulses. Supporting Information Figure S2 shows this scaling behavior using a NVIDIA GeForce GTX 1080 Ti GPU.



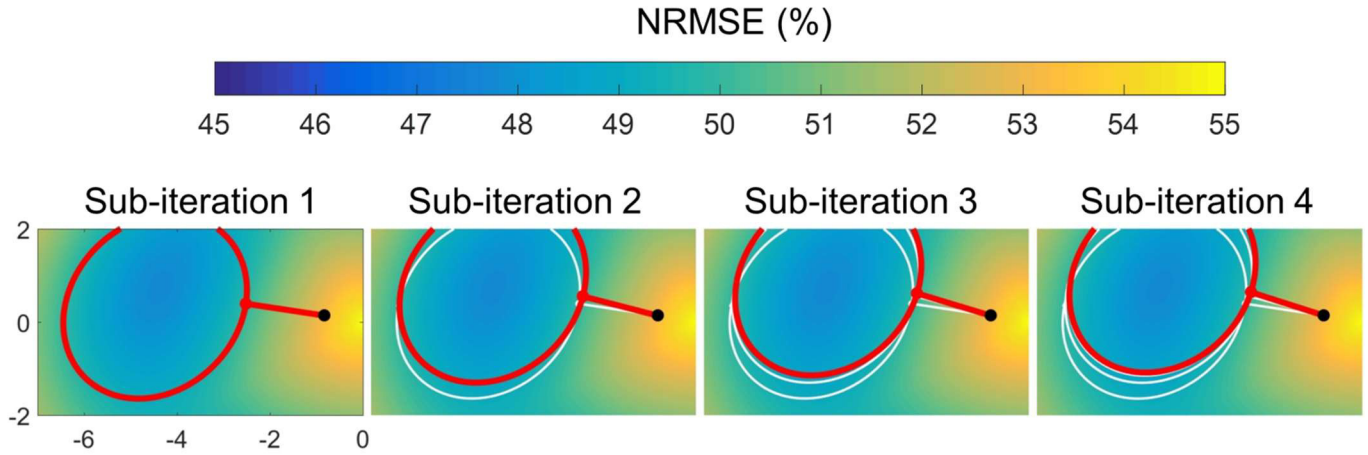


Supporting Information Figure S1: Diagram of operations in vectorized SAR oracle explained in Supporting Document B.



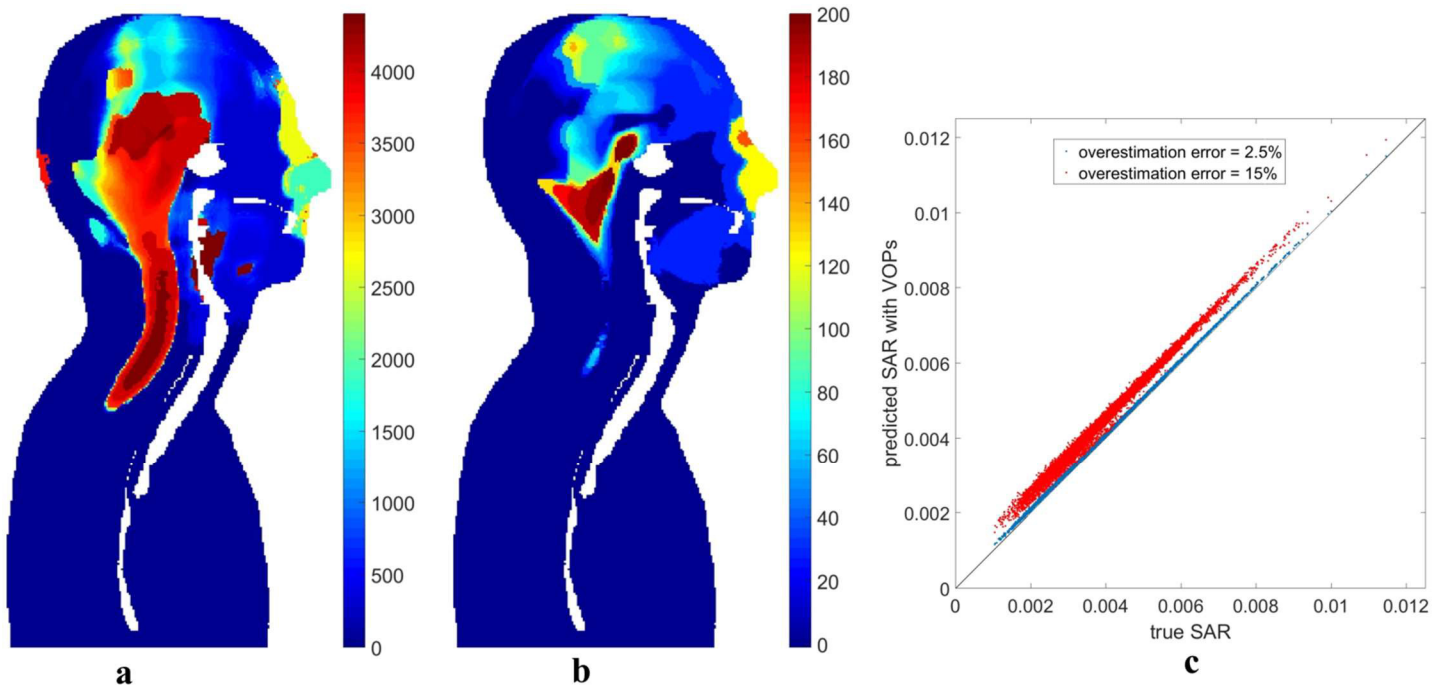
Supporting Information Figure S2: (a) Computation time using Nvidia GeForce GTX 1080 Ti for evaluation of the objective function and its gradient for different number of SAR terms and subpulses using the described vectorized oracle. The total number of SAR terms is equal to the total number of voxels in the SAR model plus the number of global SAR and per channel power terms. The total number of subpulses is equal to the total number of spokes per pulse multiplied by the total number of pulses (slices). (b) The range of pulses that can be designed with IMPULSE in approximately constant time is shown in blue based on <30 ms evaluation time in (a) since the duration of the SAR update will likely be significantly less than that of that of the FAI update when evaluation time is less than 30 ms. Accordingly, a pulse design including up to  $10^3$  to  $10^4$  subpulses and up to  $10^6$  to  $10^7$  SAR terms should be possible with IMPULSE. This corresponds, for example, to over 500 slices with 2 spokes per slice and a single high resolution ( $< 2 \text{ mm}^3$  voxel size) head model or a cluster of 25-50 lower resolution ( $> 3 \text{ mm}^3$  voxel size) models.

### C. FAI Update

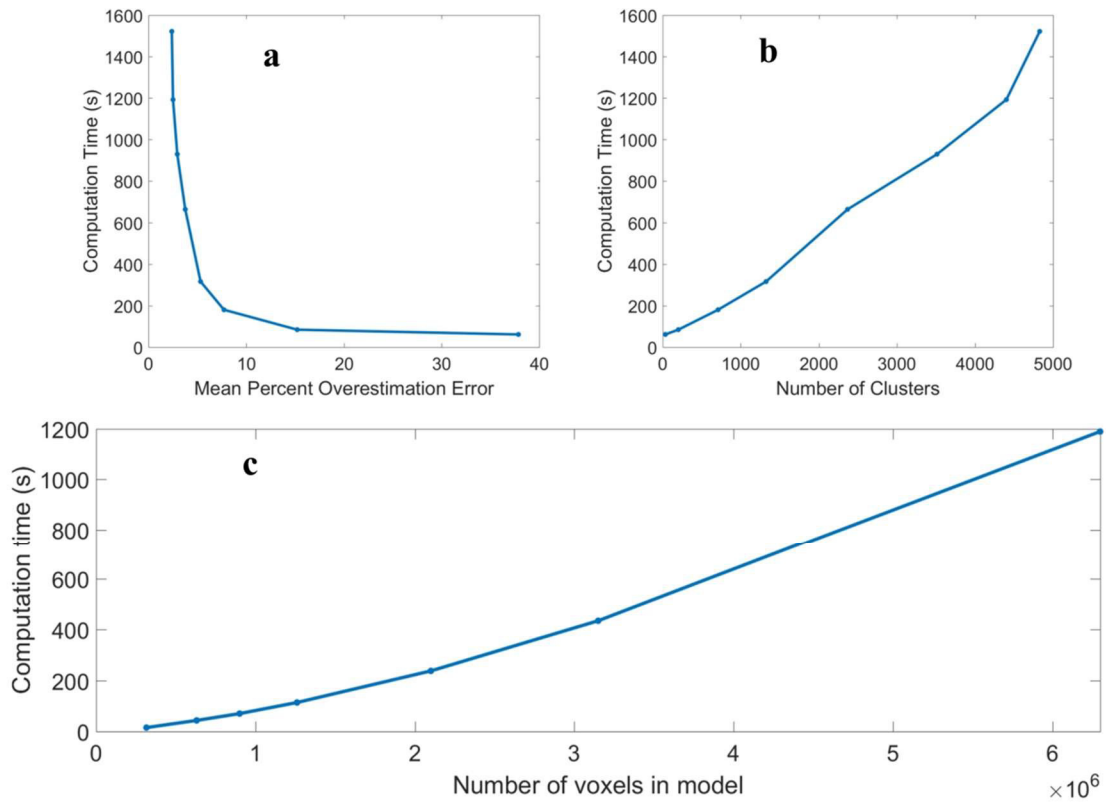


Supporting Information Figure S3: Visualization of the sub-iterations of FAI update for NRMSE tolerance of 50% for 2 channels. The background image is the NRMSE for each 2-channel weight. For a fixed set of spokes locations and target phases, the “feasible” weights that achieve the NRMSE tolerance form an ellipse (ellipse for current sub-iteration shown in red). The output of the SAR-update plus the residual ( $\underline{\mathbf{y}} + \underline{\mathbf{u}}$ ) which violates the NRMSE tolerance is shown as a black dot. The goal is to project onto the feasible set (ie. find a feasible point with minimum distance from the black dot). At each sub-iteration the black dot is projected onto ellipse defined by the current spokes locations and target phase. Then the spokes locations and target phase are updated to minimize the NRMSE with a fixed value of the channel weights leading to a new ellipse in the next sub-iteration (new ellipse shown in red, prior ellipses shown in white). The projection onto the new ellipse will have a smaller distance to the black dot than the prior ellipse. Applying this interleaved procedure of projecting and updating spokes locations and target phases gives a fast approximate solution to the FAI update despite it being nonconvex. This procedure is performed in parallel for all slices.

## D. Virtual Observation Points



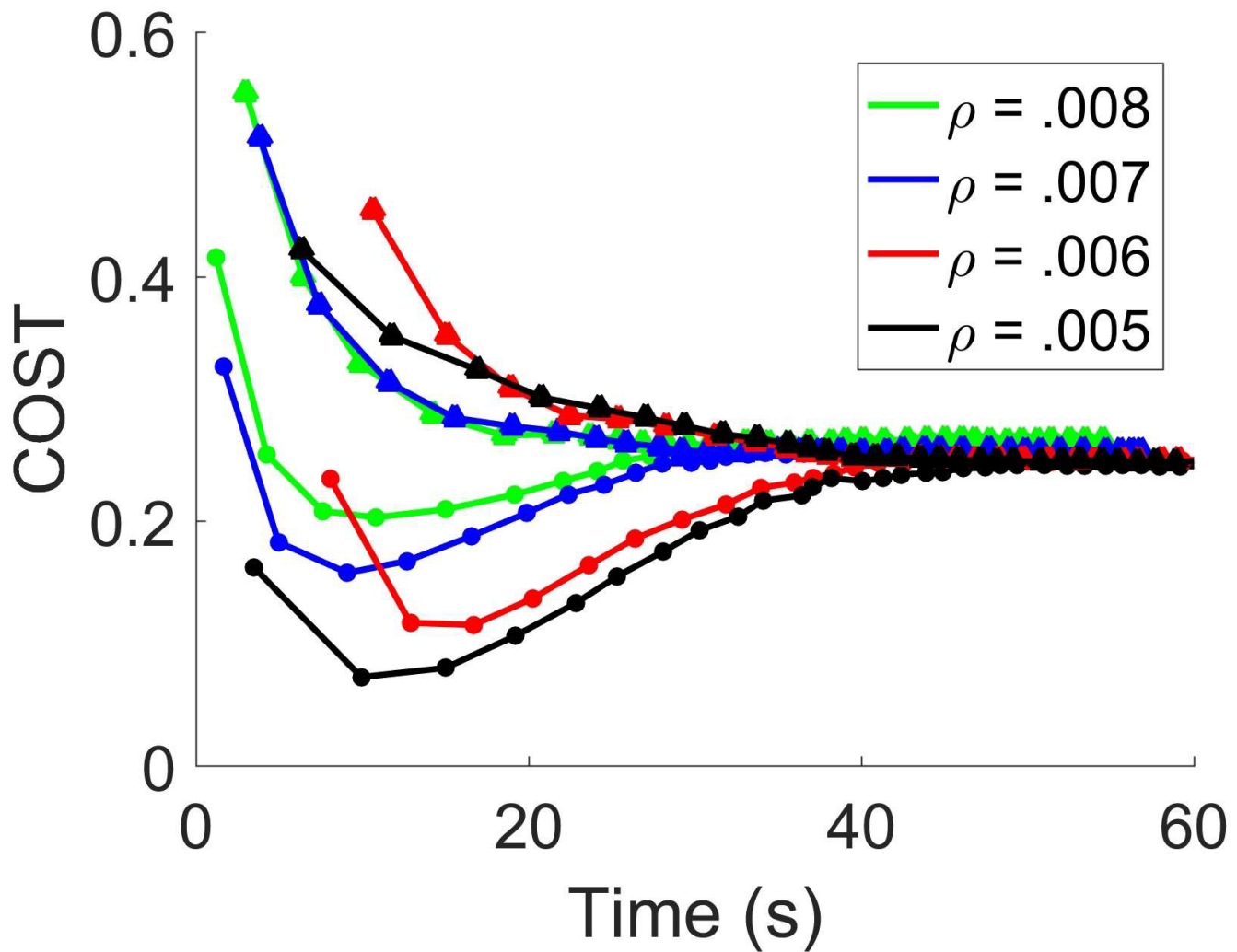
Supporting Information Figure S4 **(a)** 4401 clusters in the VOP model with 2.5% overestimation error **(b)** 200 clusters in the VOP model with 15% overestimation error **(c)** plot of true SAR vs predicted SAR with VOPs for the two VOP models for 10000 random shim weights. Original SAR model had 6371811 voxels.



Supporting Information Figure S5: **(a)** Plot of computation time for the VOP compression with 6371811 voxels for different values of overestimation error computed over 10000 random shim weights **(b)** Corresponding plot of computation time vs number of clusters for each value of the overestimation error **(c)** Plot of computation time for different numbers of starting voxels in the model with 2.5% overestimation error. Computation time increases superlinearly with number of voxels. The compression was run using C++/Eigen/OpenMP with 48 threads. The computation time significantly exceeds duration of IMPULSE for overestimation errors less than about 15% and number of voxels more than about 1 million.

## E. Tuning the IMPULSE Algorithm

The choice of the  $\rho$  hyperparameter for controlling the regularization in the SAR update will affect the rate of convergence and final value of SAR found after optimization; thus it must be chosen in a way that achieves a SAR value as close as possible to the true optimum SAR while still converging in a time frame that is suitable for on-scanner design. IMPULSE for 120 slices and 6371811 voxel Duke model was run with a range of  $\rho$  values. We chose the value of  $\rho$  for which the SAR was minimum after 45s of running IMPULSE; the duration of 45s was chosen because it is long enough to achieve convergence for a wide range of  $\rho$  while still being significantly shorter than the computation time needed for a generic optimization algorithms and being practical for on-scanner pulse design. To verify that IMPULSE solves the minSAR problem accurately, an identical problem was solved using a generic solver in MATLAB (*fmincon* using the sequential quadratic programming (SQP) algorithm taken as the ground truth) and the results were compared. As seen in Supporting Information Figure S6, for values of  $\rho$  between 0.005 and 0.008 convergence is achieved within 40-60 seconds. A value of  $\rho = 0.006$  was chosen as having lowest SAR after 45 seconds. This minimum SAR value is within 4% of the value of SAR after full convergence of IMPULSE, and within 7% of the optimum SAR found using the SQP algorithm run to convergence.



Supporting Information Figure S6: Effect of choice of  $\rho$  parameter of IMPULSE algorithm. Triangles indicate output of FAI update and dots indicate output of the SAR update. Value of  $\rho = 0.006$  seems to achieve good convergence in under 45 seconds. Larger values of  $\rho$  converge faster but to a suboptimal value while smaller values take longer to converge.

The effect of the  $\rho$  parameter is further understood through the animations in movie files rho\_small.avi, rho\_best.avi, rho\_large.avi. With a value of  $\rho$  that is too large (rho\_large.avi), the residual  $\underline{y} - \underline{z}$  (white vector) converges quickly but  $\underline{z}$  (red dot) takes longer to converge (the optimum is reached in about 8 iterations). With a value of  $\rho$  that is too small (rho\_small.avi),  $\underline{z}$  initially converges quickly but the residual takes longer to converge (the optimum is reached in about 10 iterations). With the optimal value of  $\rho$  (rho\_best.avi),  $\underline{z}$  and the residual converge at approximately the same rate (the optimum is reached in about 5 iterations).

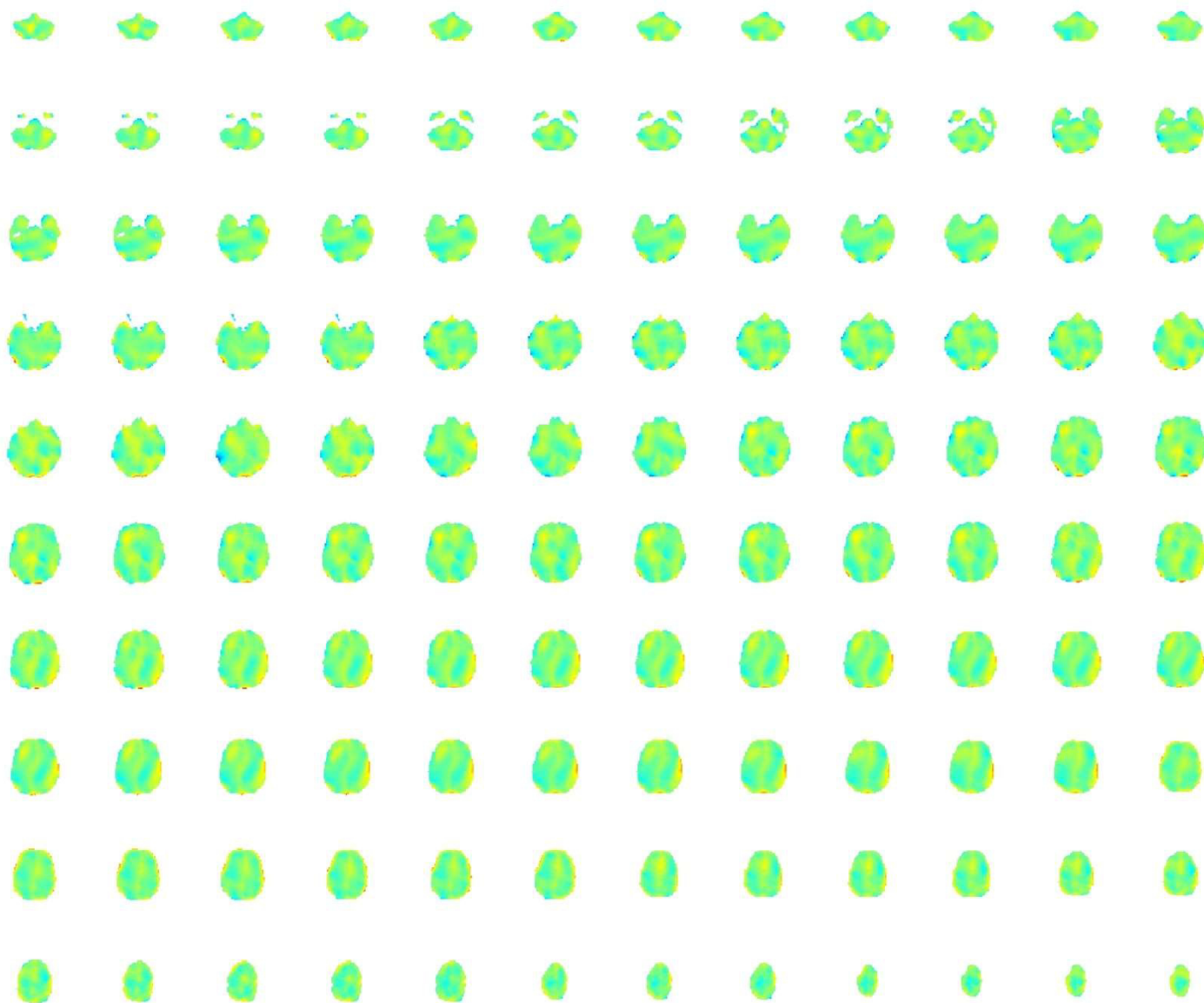




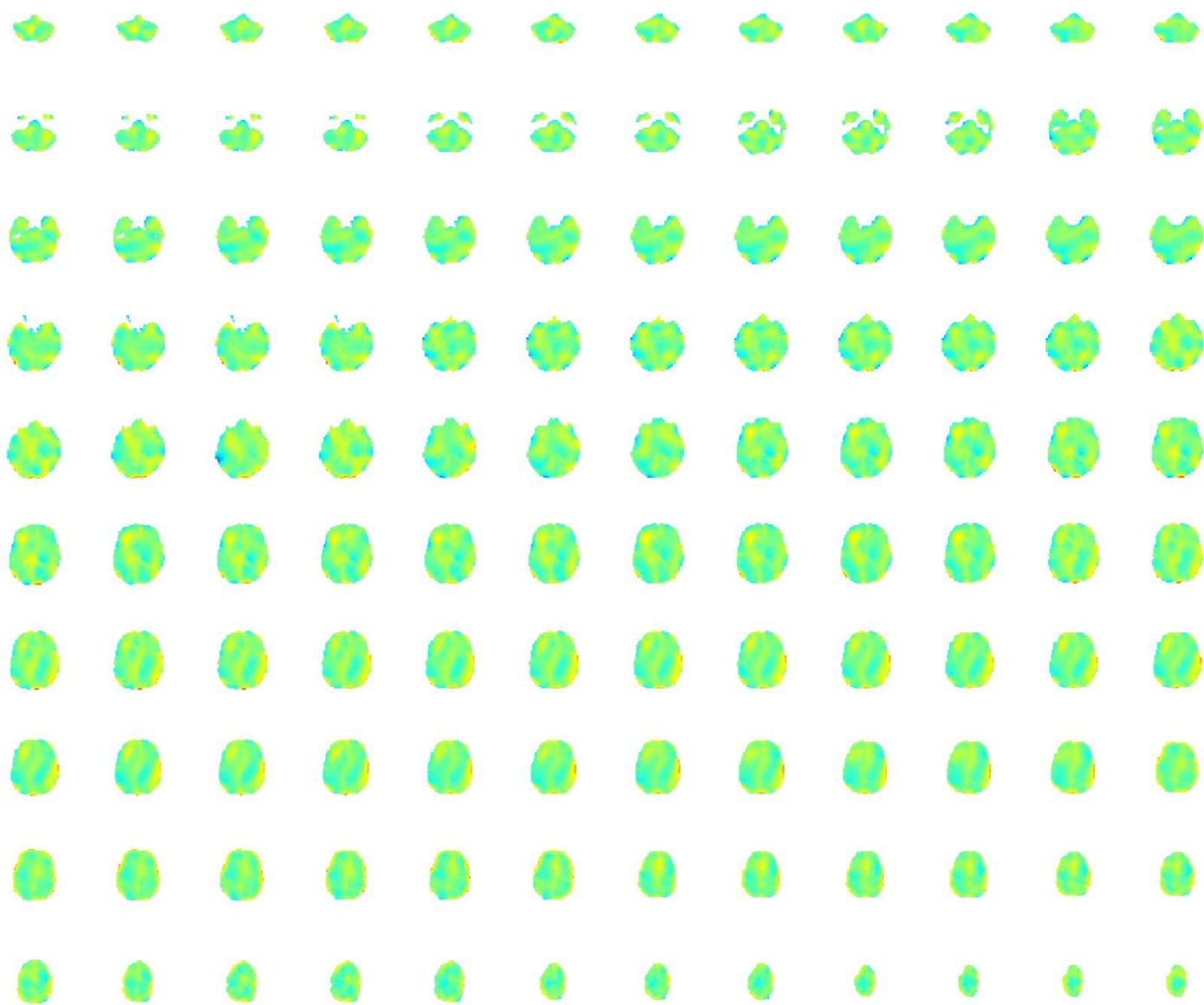
## F. Flip Angle Maps

Supporting Information Figure S7: Flip angle maps of all 120 slices for the three algorithms in Figure 7 (same colorscale).

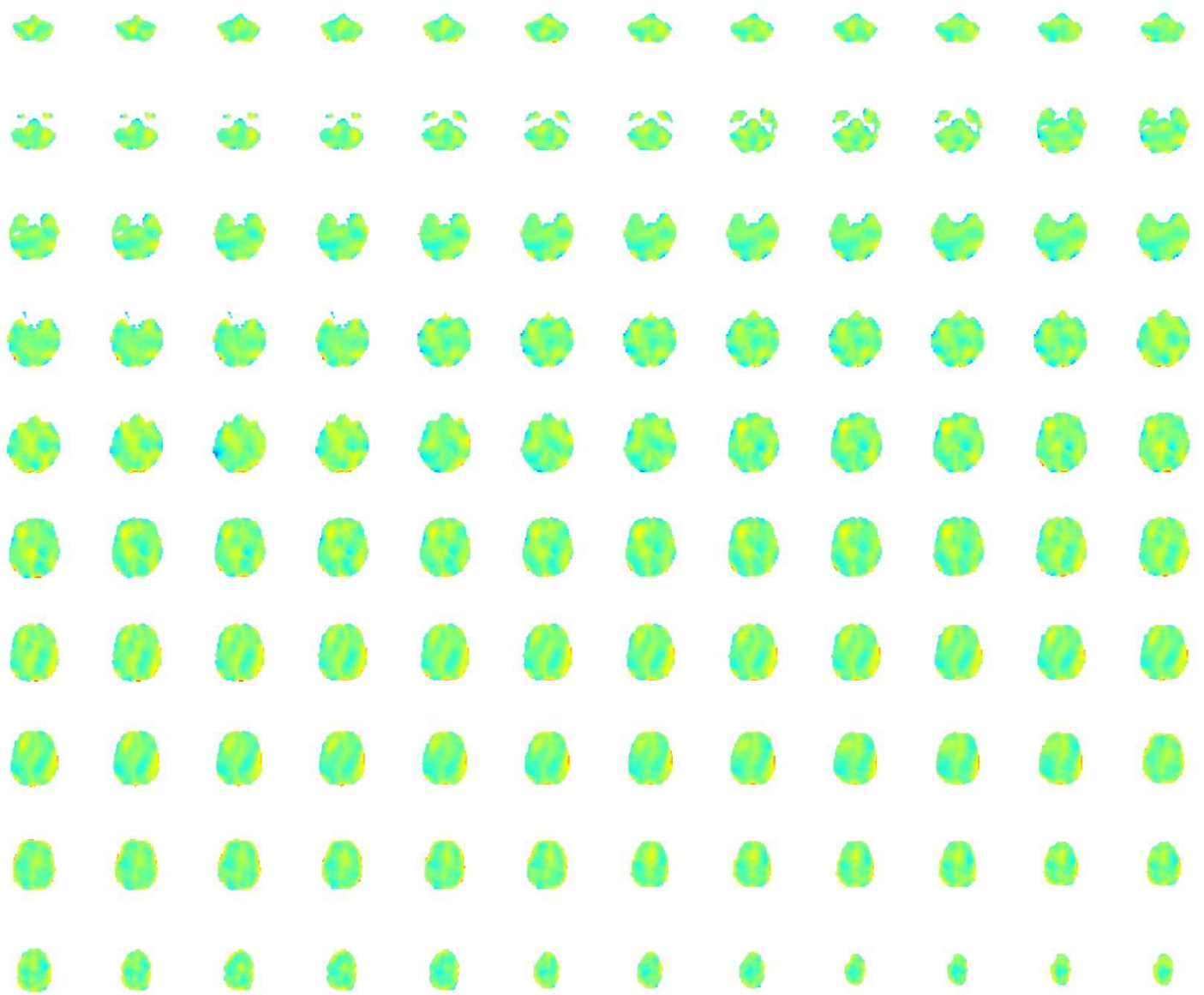
(a) Convex



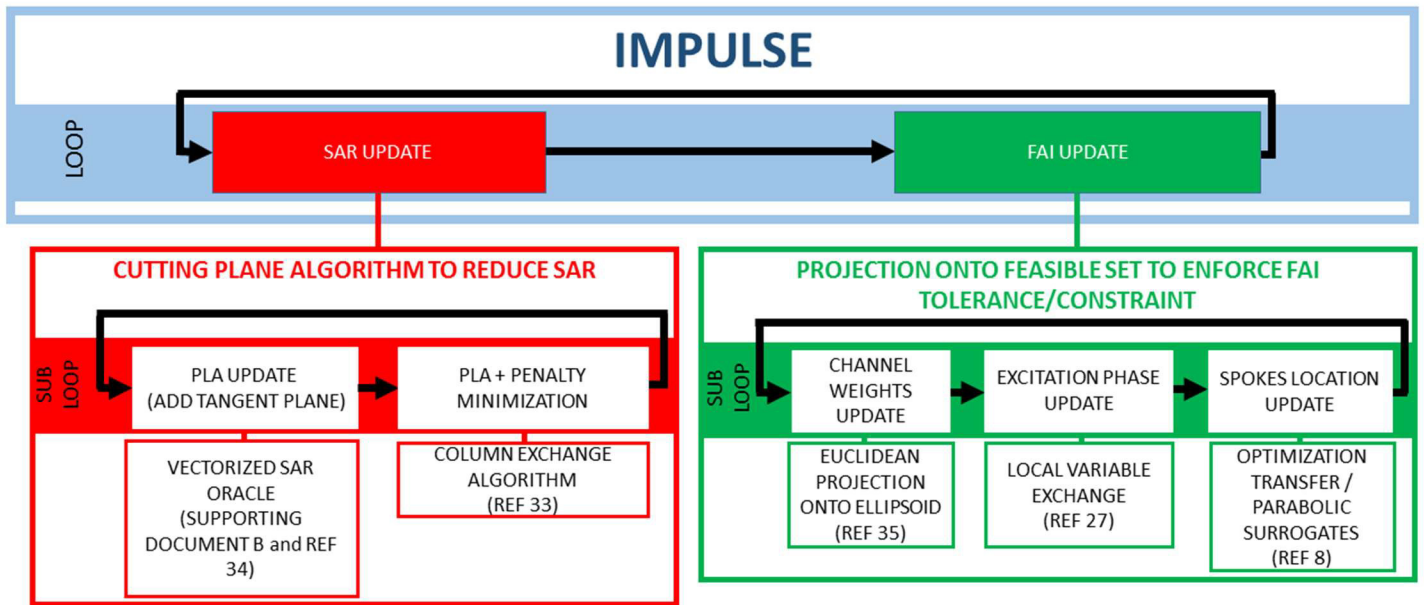
(b) SQP



(c) IMPULSE



## G. High-Level Overview



Supporting Information Figure S8: High-level overview of IMPULSE algorithm decomposed into the individual sub-algorithms and associated references for details on implementation.