

# On the optimization of imaging pipelines Schoonhoven, R.A.

## Citation

Schoonhoven, R. A. (2024, June 11). On the optimization of imaging pipelines. Retrieved from https://hdl.handle.net/1887/3762676

Version:	Publisher's Version
License:	Licence agreement concerning inclusion of doctoral thesis in the Institutional Repository of the University of Leiden
Downloaded from:	<u>https://hdl.handle.net/1887/3762676</u>

Note: To cite this publication please use the final published version (if applicable).

# **3** Auto-differentiation for CT Workflows

# 3.1 Introduction

In recent years, deep learning and other data-driven machine learning approaches have become increasingly popular in computed tomography. Deep neural networks have achieved strong results in X-ray CT applications by improving reconstruction quality [146], reducing metal artifacts [260], performing beam hardening correction [262], and classification [34, 98, 171]. The progress in deep learning has shown the power of data driven end-to-end optimization using auto-differentiation software, often in combination with hardware acceleration using graphical processing units (GPUs).

Despite promising results, deep learning approaches suffer from interpretability and reproducibility challenges. Furthermore, a serious issue is the behaviour of deep learning algorithms when presented with new data, and potential hallucination of object features [15]. These considerations have hampered the adoption of learned algorithms by experts in for example the medical domain where doctors can be reluctant to trust black-box learned approaches.

In contrast, classical (i.e. non-learned) algorithms often excel in interpretability and robustness. Additionally, classical algorithms often come with an intuition for their applicable input data domain, and the results they should produce. This can lead to expert users favouring such traditional methods over deep learning approaches. However, parameters of these algorithms are often either kept fixed or are adjusted by the user based on manual experimentation. Since these parameters are not learned in a data-driven way, they may not be optimally chosen for the particular dataset and application. In this work, we apply concepts from the philosophy that made deep learningbased methods so successful, and transfer it to CT workflows. We create workflows of computed tomography algorithms for various problems and show how those problems can be solved by end-to-end optimization based on auto-differentiation. By doing so, we reconcile classical algorithms with deep learning. We perform four case studies, representative of real-world tomography problems. To do so, we create pipelines using the following core design principles:

**End-to-end learning:** We implement the pipelines such that all pieces facilitate gradient propagation, meaning that parameters at all steps of the pipeline can be optimized jointly.

**Explicit quality criteria:** All pipelines use explicit quality criteria as objectives for optimization, thereby enabling automatic optimization of parameters.

**Declarative algorithm construction:** Each pipeline is created by building the forward model, and we optimize its parameters using auto-differentiation and generic readily-available algorithms.

Use existing building blocks: The pipelines re-use building blocks derived from both classical algorithms and deep learning methods in a way that enables gradient propagation. This allows for seamless compatibility with deep learning-based methods.

Our portfolio of case studies sketches the outline of a new generation of powerful software toolboxes that enable users to leverage the power of auto-differentiation for advanced computational CT pipeline construction.

This work is structured as follows. In Section 3.2 we explore related work in auto-differentiation software for classical methods, and approaches for learning CT algorithms in a data-driven way. In Section 3.3 we describe the methodology. In Section 3.4 we present our results in the form of four case studies. In Section 3.5 we discuss our findings and present key potential benefits that arise from our approach. We present our final conclusions in Section 3.6.

# 3.2 Related work

The idea of extending auto-differentiation techniques [72] to classical algorithms is actively investigated across several domains. In the field of robotics, for nonlinear optimization problems, Meta AI constructed Theseus [188] which is a library for building custom non-linear optimization layers that were shown to be useful for differentiable kinematics. In [36], robotic controllers are auto-tuned using the DiffTune package which works with forward-mode auto-differentiation. Furthermore, end-to-end differentiable optimization enabled the coupling of the prediction and planning module in autonomous vehicles [100]. In the field of cosmology, JAX-Cosmo [31] is a recently developed end-to-end GPU accelerated library for cosmological calculations. Using automatic differentiation, JAX-COSMO exposes derivatives for certain cosmological quantities, and enables previously impractical methods such as Hamiltonian Monte Carlo and Variational Inference. Furthermore, by embedding the cosmological algorithms in JAX, the algorithms can be run on accelerated hardware, and can benefit from automatic code optimizations and techniques such as just-in-time compilation.

Embedding classical operators in neural networks has been demonstrated as an effective technique for end-to-end learning based on auto-differentiation. In [144] the tomographic backprojection operator is embedded as an algorithm within a neural network. Here the backprojection parameters are not trainable themselves, but rather the algorithm introduces prior knowledge about image reconstruction in the neural network. From the field of seismic imaging, a 3-step reflective seismic imaging method is learned end-to-end by turning the Delay-And-Sum (DAS) operator into a network layer in [186]. Similarly, the wave-physics-formation algorithm is placed between two neural networks in [187] to facilitate single-plane seismic wave imaging (SFW).

Automatic differentiation techniques have already been applied to solve specific problems in computed tomography in recent years. In [231], beam hardening correction for X-ray microscopy of mouse bones is performed with a polynomial correction model that is optimized through a PyTorch-based differentiable FDK algorithm. An elaborate outline of algorithmic differentiation for phase retrieval is presented in [105]. In [108, 160], automatic differentiation is used in ptychography where the object wave function is obtained with a gradient-based method by minimizing the ptychography loss for each pixel. In [55] the 3D reconstruction problem for objects beyond depth-of-focus (DOF) is formulated as a minimization problem with a data fidelity and total variation term, which is solved with a gradient descent algorithm. How automatic differentiation techniques can be used for different imaging modalities is shown in [77], where compressive sensing, single image super resolution (SISR), and ptychography reconstructions are covered. For tomography, the reconstruction is obtained with a gradient based approach by minimizing a total variation functional, and the authors show how this is beneficial for sparse and limited angle data.

In the field of nanotomography, the auto-differentiation framework Adorym [54] for flexible reconstruction has been developed. In Adorym, a flexible forward model allows for optimization of experimental parameters such as probe position, object tilt, absorption/refraction relation coefficient, and propagation distance. The authors show improved reconstruction quality for ptychography and multi-distance holography reconstructions, by finetuning these experimental parameters. For conventional CT, they accelerate ART with gradient descent optimizers and acquire improved results over FBP.

Compared to the related work outlined above, the scope and focus of the present paper is more general. We focus on the core design principles underlying CT workflows using auto-differentiation and demonstrate the efficacy of the approach through a varied set of four case studies.

# 3.3 Methodology

### 3.3.1 Auto-differentiation

For the case studies included in this paper, we implement the CT workflows in the auto-differentiation framework PyTorch [181]. An auto-differentiation system breaks down a program in a series of primitive operations for which fixed procedures are known to compute derivatives. The series of primitive computations is collected in a computational graph. Most auto-differentiation systems, including PyTorch, trace the computational graph implicitly during the forward computation through the program. After the forward computation, the error or loss is computed. Here, we use gradient-based optimization, meaning that an auto-differentiation package needs to obtain partial derivatives of the loss with respect to the parameters.

Often each primitive function in an auto-differentiation package specifies vector-Jacobian products. Suppose two quantities are related by a primitive function  $f(\mathbf{x}) = \mathbf{y}$ , and  $\mathbf{y}$  is used further in the computation. Denote by  $\bar{\mathbf{y}}$  the derivative of a loss L with respect to  $\mathbf{y}$ . A vector-Jacobian product defines a way to express the derivatives of L with respect to  $\mathbf{x}$ , and is defined as

$$\frac{\partial L}{\partial x_j} = \sum_i \frac{\partial y_i}{\partial x_j} \frac{\partial L}{\partial y_i}, \qquad \bar{\mathbf{x}} = J^T \bar{\mathbf{y}},$$

where J is the Jacobian. For a primitive operation f, the gradient of the input  $\bar{\mathbf{x}}$  can be calculated from the output gradient  $\bar{\mathbf{y}}$ , the input  $\mathbf{x}$  and the output  $\mathbf{y}$ . For example, in the case of y = f(x) = -x the gradient of the input can be computed as  $-\bar{y}$ , for  $y = f(x) = e^x$  it will be  $y \cdot \bar{y}$ , and in the case of  $y = \log(x)$  it will be  $\frac{\bar{y}}{x}$ . The procedure of computing gradients in reverse (starting from  $\bar{L} = 1$ ) is called back propagation. For gradient descent, parameters are updated by

$$\mathbf{x} \leftarrow \mathbf{x} - \lambda \bar{\mathbf{x}}$$

for some step size  $\lambda$  (often called the learning rate in deep learning). Variants of gradient descent exist such as Nesterov's accelerated gradient descent [162]. In our experiments we will either use (stochastic) gradient descent, or Adam [113].

A drawback of auto-differentiation frameworks can be excessive memory consumption because copies of many intermediate outputs are stored for fast computation of the back propagation algorithm. This problem can be addressed by several general approaches. Gradient checkpointing stores only a subset of the intermediate outputs for gradient computation [35]. Another approach involves just-in-time (JIT) compilation where a compiler attempts to optimize the computational graph into a more memory and compute efficient set of instructions. In this work, we apply PyTorch implementations of both techniques for certain memory or computationally expensive operations.

## 3.3.2 Computed tomography

In computed tomography [106] a 3D image of an object is recovered from a series of projections that are taken at different angles. Tomographic reconstruction can be modelled as the problem to recover an object volume  $\mathbf{x} \in \mathcal{X} := \mathbb{R}^{N_x \times N_y \times N_z}$ from the measured projection data  $\mathbf{y} \in \mathbb{R}^{N_\theta \times N_u \times N_v}$ . Here,  $N_u$  and  $N_v$  are the number of detector rows and columns, and  $N_\theta$  is the number of projection angles. The projection process can be approximated by a linear operator A, and can be expressed as a matrix using the aforementioned discretization

$$A\mathbf{x} = \mathbf{y},\tag{3.1}$$

where  $\mathbf{x}$  and  $\mathbf{y}$  are collapsed to a vector. For parallel beam tomography the object can be reconstructed by the commonly used filtered backprojection algorithm (FBP)

$$\mathbf{x}_{\text{FBP}} = A^T (\mathbf{h} * \mathbf{y}). \tag{3.2}$$

Here, the projection data is convolved with a 1D filter  $\mathbf{h} \in \mathbb{R}^{N_v}$  (Ram-Lak in this study), and subsequently backprojected by applying  $A^T$ . For circular cone-beam tomography, the object can be recovered by the Feldkamp-Davis-Kress (FDK) [59] algorithm, where the projection data is weighted in order to compensate for the diminishing intensity at distance from the detector center. An alternative to direct methods are iterative methods that reformulate the equation system 3.1 as an optimization problem of the form

$$\mathbf{x}^* = \underset{\mathbf{x}\in\mathcal{X}}{\arg\min} \|A\mathbf{x} - \mathbf{y}\|_2^2.$$
(3.3)

Variational methods additionally aim to incorporate prior knowledge in the form of a regularization term  $R(\mathbf{x})$  to the functional, e.g.,

$$\mathbf{x}^* = \underset{\mathbf{x}\in\mathcal{X}}{\arg\min} \|A\mathbf{x} - \mathbf{y}\|_2^2 + R(\mathbf{x}).$$
(3.4)

## 3.3.3 CT workflows

The CT workflows that we consider here contain input data in the form of projection images, contain several data processing steps of which at least one is a reconstruction step, and contain an objective function that scores the final result. We consider potentially non-sequential workflows, i.e., the data processing graph can contain several branches, or cyclical sections. The computational blocks can contain parameters that we want to update in order to minimize the objective function (see Figure 3.1). Those parameters must be real or complex numbers in order to facilitate gradient-based optimization.

For the four case studies, we have implemented the computational blocks to be differentiable with respect to the learnable parameters. Non-differentiable steps can potentially be replaced by a differentiable approximation function. For example, to make thresholding (and segmentation) differentiable, we implement the thresholding operation using a hyperbolic tangent

$$\tau_{\gamma}(\mathbf{x},t) = \frac{1}{2} \left( 1 + \tanh\left(\gamma\left(\frac{\mathbf{x}}{t} - 1\right)\right) \right), \qquad (3.5)$$

where the volume  $\mathbf{x}$ , and threshold t have been scaled to [0, 1]. The parameter  $\gamma$  regulates the sharpness of the clipping.



Figure 3.1: Example CT workflow diagram.

## 3.3.4 Software implementation

To compute gradients end-to-end for workflows that contain tomographic (back)projection operators, we will make use of the matrix identity  $\nabla A \mathbf{x} = A^T \nabla \mathbf{x}$ . Tomographic projection operations are implemented in the ASTRA toolbox [1] in a computationally efficient GPU accelerated manner. The tomosipo package [90] implements PyTorch support for ASTRA, and contains projection operators that propagate gradients in PyTorch using the aforementioned identity. Here, we use tomosipo projection operators in our workflows to propagate gradients end to end. In addition to reconstruction algorithms, the workflows we construct in this work contain other classical CT algorithms, such as Paganin's phase retrieval, phase contrast projection, and spectral projectors. We implemented these algorithms in PyTorch which facilitates auto-differentiation, and makes them GPU compatible, and the code is made publicly available for all case studies [216].

# 3.4 Case studies

## 3.4.1 Rotation axis alignment

#### Introduction

In the first case study, we consider the problem of alignment in tomographic reconstruction [123]. Various experimental factors can introduce errors in the geometric parameters that are used to perform the reconstruction. A common occurrence of this effect in CT is a misalignment of the rotation axis position. In

the case of a rotation axis misalignment  $\delta$ , for a ray parametrized by angle  $\omega$ and position vector a from scanning set  $\Gamma$ , the projection p of the object function  $f: \mathbb{R}^3 \to \mathbb{R}$  is given by the Radon transform as

$$p(\boldsymbol{a},\boldsymbol{\omega}) = \int_0^\infty f(\boldsymbol{a} + \boldsymbol{\delta} + t\boldsymbol{\omega}) \, dt, \quad \boldsymbol{\omega} \in \mathbb{S}^2, \quad \boldsymbol{a} \in \Gamma.$$
(3.6)

This discrepancy (most prominently a lateral shift) between the assumed and the true rotation axis position introduces severe artifacts in the resulting reconstruction, which appear differently depending on the acquisition geometry. In cone-beam CT this leads to blurring and "doubling" of object features, which significantly affect the sharpness of the reconstruction.



(d) Contrast measure during optimization (e) Rotation axis shifts during optimization

Figure 3.2: (a) Shepp-Logan phantom and its FDK reconstructions obtained from projection data simulated using (b) shifted rotation axis and (c) rotation axis position compensated by the proposed gradient-based optimization method. (d) Image variance-based contrast measure and (e) estimated rotation axis shift during the optimization.

#### Experiments

Here we propose to optimize for the rotation axis position by minimizing the magnitude of misalignment-induced artefacts as measured by an appropriate measure of reconstruction quality [50] [76]. In cone-beam CT, to account for edge blurring and "doubling" artifacts, we use a well-known image contrast measure based on voxel variance [73]. The intuition behind this metric is that a sharply focused reconstruction has a large spread in the intensity histogram, whereas upon blurring intensities concentrate around the same gray value. To enable gradient-based optimization for this task, we adjust the lateral shift of the assumed rotation axis position by shifting all projection images in the opposite direction using a bicubic interpolation-based image shift operator implemented in PyTorch. The shifted projections are then backprojected and the quality of this intermediate result is quantified. Since all the pieces of the described workflow are implemented in a differentiable manner within PyTorch, a generic gradient-based optimization algorithm (SGD) can now be used to find the shift of the rotation axis that maximizes the reconstruction quality. A diagram of the proposed pipeline is shown in Figure 3.3. Both experiments in this section were performed on a workstation with 64GB RAM, an NVidia GeForce GTX 1070 GPU, and Intel i7-7700K CPU.



Figure 3.3: CT workflow for self-supervised rotation axis alignment.

#### Simulated experiment: Shepp-Logan phantom

To test the proposed gradient-based rotation axis alignment method, we first employ simulated data based on a  $512 \times 512$  px Shepp-Logan phantom (Figure 3.2(a)) that is projected using a geometry with a rotation axis shift of 3 px. As can be observed in Figure 3.2(b), even a minor misalignment in the axis position introduces significant blurring in the reconstruction. The proposed gradient-based contrast optimization method successfully compensates the rotation axis shift (Figure 3.2(c-e)), converging in only about 3 iterations. The proposed method is computationally efficient as running 6 iterations took 97 ms on our system.

#### Real-world experiment: High-resolution cone-beam CT of a walnut

Next, we evaluate the performance of the proposed rotation axis shift compensation method on experimental data using an open dataset of high-resolution cone-beam CT of walnuts acquired at the Flex-ray lab [39, 119]. Figure 3.4(a) demonstrates an FDK reconstruction obtained with the geometry parameters specified in the dataset metadata, and Figure 3.4(b) shows the reconstruction after axis alignment method

Voxel variance



(a) FDK (before optimization)





Training iteration

(c) Contrast measure during

optimization

Training iteration (d) Rotation axis shifts during optimization

3

Figure 3.4: (a) FDK reconstruction of the walnut dataset. (b) FDK reconstruction after rotation axis position was compensated by the proposed gradient-based optimization method. (d) Image variance-based contrast measure and (e) estimated rotation axis shift during the optimization.

has been applied. It can be observed that although there are no obvious artifacts present in the initial reconstruction (which makes this kind of misalignment easy to miss with manual inspection), the rotation axis position optimization significantly improves the resolution of the reconstruction. The improved resolution brings up details that were not visible before, such as pores in the central part of the walnut that can be seen in the zoom-in of Figure 3.4(b). The walnut slice is  $550 \times 550$  px, with 501 projection angles, and running 6 iterations took 295 ms on our system.

To summarize, our workflow for rotation axis position alignment, implemented in an automatic differentiation framework, results in an intuitive formulation of the axis alignment problem. In addition, the resulting method enjoys a quick convergence, can retrieve sub-pixel axis shifts in a straightforward manner, and can be easily extended to multivariate optimization if other geometry parameters, such as rotation axis tilt and the cone angle of illumination, are included in the workflow.

## 3.4.2 Phase retrieval

#### Introduction

In the second case study we apply an end-to-end optimization approach to phase contrast imaging (PCI) [57]. PCI can reach nanometric resolution in tomographic imaging [251], and requires an additional reconstruction step known as phase retrieval. In PCI, the image is reconstructed based on changes to the wave front due to the material that is present along the wave path. A widely used experimental PCI setup is phase propagation-based imaging where projections are acquired from several different distances using a coherent beam. Next, a phase retrieval algorithm is used to calculate phase maps.

Here, we will consider Paganin's phase retrieval algorithm [172], which requires the material refractive index  $\delta$ , and attenuation  $\beta$  to be known. Paganin assumes a single material object, and retrieves the projected thickness of the object **T** by

$$\mathbf{T}(\mathbf{r}_{\perp}) = -\frac{1}{\beta} \log \left( \mathcal{F}^{-1} \left( \frac{\beta \mathcal{F}\{I(\mathbf{r}_{\perp}, z = R_2)\}/I_0}{R_2 \delta \|\mathbf{k}_{\perp}\| + \beta} \right) \right).$$
(3.7)

Here I is the intensity function,  $I_0$  the incident intensity,  $\mathbf{r}_{\perp}$  the position vector perpendicular to the optical axis,  $\mathbf{k}_{\perp}$  the wave vector, and  $R_2$  the source-detector distance. A common practice is to divide both the numerator and the denominator inside the inverse Fourier transform of equation 3.7 by  $\beta$ . The resulting fraction  $\delta/\beta$  is sometimes designated as  $\alpha$ . Expert users will often pick  $\alpha$  manually to get a good reconstruction. However, this can be a time-consuming process, it is subjective, and it can make it more difficult for other researchers to reproduce results. Therefore, in this section we will show that we can optimize both  $\beta$  and  $\delta$ in an unsupervised way for a clear objective. We do so by constructing a pipeline of operators that propagate gradients end-to-end.

#### **Experiments**



Figure 3.5: PCI pipeline for self-supervised optimization of  $\beta$  and  $\delta$ .

In the experiments, we use a pipeline that takes raw projections as input, and performs Paganin phase retrieval to produce phase maps. A reconstruction based on the phase maps is made with filtered backprojection (FBP), and a binary segmentation is made with an implementation of Otsu's method [170] that makes use of equation 3.5. We use binary segmentation since Paganin assumes a single material.

After segmentation, using the same refraction and attenuation indices, projections are simulated based on the segmentation using a wave propagation projector. Finally, the simulated projections are scaled so that their mean and standard deviation align with the raw input. As a loss function we take the mean-squared error loss between the scaled simulated projections and the original input projections. A diagram of the pipeline is given in Figure 3.5. Both experiments in this section were performed on a server with 384GB RAM, an NVidia Titan RTX GPU, and two Intel Xeon Gold 6130 CPUs.



Figure 3.6: (a) Slice through calcium carbonate simulated phantom. (b): FBP reconstruction using correct  $\beta$  and  $\delta$ . (c): FBP of phase projections retrieved with the material indices for water that were used as initialization. (d): FBP of phase projections retrieved with learned material parameters after optimization.



Figure 3.7: (Left): L2-loss between input phase maps and output simulated phase maps per iteration. Attenuation  $\beta$  (middle), and refraction  $\delta$  (right) value per iteration.

#### Simulated experiment: Calcium carbonate cube

First, we perform a simulated experiment on a 3D phantom (192<sup>3</sup> volume, 572 angles) of a calcium carbonate hollow cube, with a smaller cube attached to one of its sides (see Figure 3.6(a)). Next, we simulate phase contrast images using the material parameters of calcium carbonate, and add Poisson noise. As reference, we show the FBP reconstruction of the phase maps acquired with Paganin with the correct  $\beta$  and  $\delta$  in Figure 3.6(b). We initialize the attenuation and refraction indices to those of water. We use a gradient descent algorithm (Adam) optimization algorithm to update  $\beta$  and  $\delta$ . For a smoother optimization we optimize these parameters on an exponential scale.

The FBP reconstructions before and after optimization are given in Figure 3.6. We see that the initial FBP reconstruction using water material indices has a halo artifact. This artifact is no longer visible after optimization. The optimization loss and the attenuation and refraction values (logarithm) per iteration are given in Figure 3.7. While the reconstruction quality has improved, the pipeline is not able to fully recover the original material index values; the final values after optimization



Before

After

Figure 3.8: Hydrogen fuel cell [45] reconstructed with FBP with default material parameters for water (before), and after optimizing parameters end-to-end (after).

are  $\ln(\beta) = -21.50$  and  $\ln(\delta) = -15.69$ , whereas the original values for calcium carbonate are  $\ln(\beta) = -19.59$ ,  $\ln(\delta) = -13.94$ . However, the value of  $\alpha = \delta/\beta$  is only 17.6% off the original. Running 120 iterations of optimization on this pipeline took 16 minutes and 48 seconds on our system.

#### Real-world experiment: Hydrogen fuel cell

We perform a real-world data experiment on a hydrogen fuel cell dataset acquired at the TOMCAT beamline of the Swiss Light Source (PaulScherrer Institut)[45]. The experiment is performed on a central slab of the full volume of size  $160 \times 1001 \times 1476$ to reduce computation time. For the experiment we use the same setup as for the simulated experiment, and again initialize the attenuation and refraction indices to those of water. In Figure 3.8 we show reconstructed central slices of the fuel cell before and after optimization. We see that after optimization small-scale features such as bubbles are visible. The zoomed images show that the initial reconstruction is blurred, and the final reconstruction after optimization is much sharper. In Figure 3.9 we display the optimization loss and the attenuation and refraction values per iteration. Running 200 iterations of optimization on this pipeline took 1 hour and 38 minutes on our system.

Overall the experiments show that we can use a generic gradient based approach to optimize for reconstruction quality in an self-supervised way in phase contrast imaging. We showed on both a simulated and real-word dataset that the pipeline can be optimized for a clear objective, as opposed to manual selecting the  $\alpha = \delta/\beta$ parameter. First, this reduces the work for operators as they no longer have to



Figure 3.9: (Left): L2-loss between input phase maps and output simulated phase maps per iteration. Attenuation  $\beta$  (middle), and refraction  $\delta$  (right) value per iteration.

tune parameters manually. Second, this makes the procedure more reproducible for different datasets.

## 3.4.3 Beam hardening correction

#### Introduction

For the third case study, we implemented self-supervised correction of artifacts introduced by beam hardening [23, 91]. The number of materials is denoted by N, with attenuation coefficients  $\mu_n(E)$ , and the beam spectrum has  $E_m$  energy bins with intensities  $I_e$ . Furthermore, let  $l_{i,j}$  be the intersection length of ray i $(i = 1 \dots D)$  with voxel j  $(j = 1 \dots J)$ ,  $d_j$  the relative density of voxel j, and  $s_{n,j}$ a variable that is 1 if voxel j contains material n, and 0 otherwise. Then we can denote for a ray i and a material n the projected relative density  $P_{i,n}$ , and the monochromatic measured intensity by Beer-Lambert as

$$P_{i,n} = \sum_{j=1}^{J} l_{i,j} d_j s_{n,j}, I_{\text{mono},i} = I_0 e^{-\sum_{n=1}^{N} \mu_n(E_0) P_{i,n}},$$
(3.8)

for a monochromatic beam with energy  $E_0$  and intensity  $I_0$ . The measured polychromatic intensity (discrete) is given by

$$I_{\text{poly},i} = \sum_{e=1}^{E_m} I_e e^{-\sum_{n=1}^N \mu_{n,e} P_{i,n}},$$
(3.9)

for a polychromatic beam with energies e and intensities  $I_e$ .

For a monochromatic X-ray, the attenuation coefficient  $\mu$  is linearly related to the thickness of the object by Beer-Lambert's law. However, in practice Xray beams are often polychromatic and this relation is no longer linear as lower energy photons get absorbed more than higher energy photons. Therefore, as a polychromatic beam travels through an object it "hardens", i.e., the average photon energy increases. If this effect is not taken into account by the reconstruction algorithm, it causes streaking and cupping artifacts because the lower absorbance due to higher average energy is mistakenly reconstructed as a lower density material.

#### Experiments

The aim of the case study is to highlight how combining an explicit forward model with a few lines of code and generic gradient-optimization can result in sophisticated algorithm construction, even for cases that used to require lengthy, hand-crafted implementations. In the experiments we will perform unsupervised beam hardening correction by learning the beam spectrum, and the energy-dependent attenuation coefficients per material. We will base the study on [66], where three different iterative unsupervised beam hardening correction algorithms are proposed and compared. We use the best performing algorithm in [66] as a comparison; the *iterative sinogram preprocessing* (ISP) method. The ISP method is an iterative scheme with multiple steps, such as a bruteforce segmentation step, a local optimization step, and a reconstruction update step. For a detailed explanation of ISP we refer to [66]. ISP assumes that the number of materials is known beforehand, but no knowledge of the beam spectrum, or attenuation coefficients is assumed.



Figure 3.10: CT pipeline for self-supervised beam hardening correction.

As opposed to constructing a specialized algorithm such as ISP, we compare to a generic gradient-based approach where beam spectrum, attenuation coefficients, and thresholds are learned jointly. The generic approach is self-supervised using the same loss function as ISP

$$\varphi(\mu, \mathbf{I}, \mathbf{s}, \mathbf{d}) = \frac{1}{D} \sum_{i=1}^{D} \left( \log\left(\frac{I_{\text{poly},i}^{\text{meas}}}{I_0}\right) - \log\left(I_{\text{poly},i}^{\text{sim}}\right) \right)^2, \quad (3.10)$$

which is the  $L_2$ -loss between the simulated polychromatic projections and the original input projections. To initialize the thresholds we use Otsu's method on an initial reconstruction  $R_0$  made with FBP or FDK. After the initialization we iterate as

1. Update segmentation thresholds, attenuation coefficients, and beam intensities,

$$\mathbf{s}^{k}, \mu^{k}, \mathbf{I}^{k} = \underset{\mathbf{s}, \mu, \mathbf{I}}{\operatorname{arg\,min}} \ \varphi(\mu, \mathbf{I}, g(R^{k-1}, \mathbf{s}), \mathbf{d} = 1),$$
(3.11)

where  $g(R^{k-1}, \mathbf{s})$  is the differentiable thresholding function described before (equation 3.5),  $\varphi$  is the ISP loss function, and arg min is performed with gradient descent.

2. Update the sinograms and corrected reconstruction  $\mathbb{R}^k$  as for ISP.

A diagram of the workflow is shown in Figure 3.10. By using an implementation of the spectral projector that can propagate gradients, we can optimize the attenuation, intensity, and thresholds jointly. This improves the complexity of the algorithm as the thresholding step in ISP is exponential in the number of materials and computationally expensive. For our experiments, we created a PyTorch based implementation of ISP and our generic gradient-based approach. Both experiments in this section were performed on a workstation with 64GB RAM, an NVidia RTX 2070 Super GPU, and AMD Ryzen 7 3800X CPU.



Figure 3.11: (Left): Barbapapa phantom consisting of PMMA filled with aluminium cylinders. (Right): FBP reconstructions of simulated spectral projections (with added Gaussian noise).

#### Simulated experiment

We perform a simulated experiment based on the Barbapapa phantom [66] where we created a simulated phantom (see Figure 3.11) of polymethyl methacrylate (PMMA) filled with aluminium rods with size  $256 \times 256$ . Next, we use a spectral projector with a simulated effective beam spectrum to simulate beam hardening artifacts. Finally, we add 2% Gaussian noise on the projections (see Figure 3.11 for reconstruction of noisy projections). To perform beam hardening correction we performed 40 steps of ISP. We ran the generic gradient-based approach for the same amount of objective function evaluations. In Figure 3.12 we show the FBP reconstructions of the corrected sinograms for both algorithms, and the accompanying material segmentations. We see that both methods reduced cupping and streaking artifacts. Furthermore, Figures 3.12(d) to 3.12(f) show a significant improvement in segmentation quality for both methods, which is close to the original phantom.

Line profiles for each of the three reconstructions are given in Figure 3.13(a). The line profiles confirm that the cupping artifacts are significantly reduced for both the PMMA material and the aluminium rods. In Figure 3.13(b) we plot the self-supervised loss for both methods. Both methods reach a similar optimum, but the combined gradient method shows slightly faster convergence. Note that the combined gradient method also runs twice as fast, which is due to the missing brute-force threshold selection step. This is also reflected in the runtimes; ISP ran for 16 minutes and 15 seconds, whereas the generic gradient-based approach ran for 8 minutes and 9 seconds.

#### Play-Doh foreign object X-ray CT dataset

We also evaluate both beam hardening correction methods on a real-world X-ray CT dataset of Play-Doh objects filled with stones [257, 258]. The experiment was



(d) Original FBP segmentation (e) ISP segmentation

(f) CG segmentation

Figure 3.12: (Top): Barbapapa phantom reconstructed with the Filtered BackProjection (FBP) algorithm for ISP and a combined gradient-based (CG) approach. (Bottom): Material segmentations of corrected reconstructions.



Figure 3.13: (a) Three line profiles through the original, ISP corrected, and combined gradient corrected Barbapapa reconstructions. (b) Optimization losses per  $\varphi$  evaluation.

performed on a  $478 \times 478$  central slice. The Play-Doh objects exhibit significant cupping artifacts. In Figure 3.14 we show the FBP reconstructions of the corrected



Figure 3.14: Play-Doh and stones reconstructed with the Filtered BackProjection (FBP) algorithm for ISP and a combined gradient-based (CG) approach.

sinograms for both algorithms. We see that both methods reduced cupping, which is also clearly visible on the line profiles shown in Figure 3.15(a). In Figure 3.15(b) we plot the self-supervised loss for both methods. Both methods reach a similar optimum, but the combined gradient method shows a more stable convergence. The jumps in the loss curve for ISP happen when a new optimal set of thresholds is determined with brute-force calculation after the local minimization steps. Since the combined gradient method jointly optimizes all parameters every step, it does not suffer from such jumps. The improved computational efficiency is more prominent due to faster convergence of the local optimization step; ISP ran for 19 minutes and 51 seconds, whereas the generic gradient-based approach ran for 4 minutes and 35 seconds.



Figure 3.15: (a) Three line profiles through the original, ISP corrected, and combined gradient corrected Play-Doh reconstructions. (b) Optimization losses per  $\varphi$  evaluation.

Overall the experiments show that by using a generic gradient-based approach for all parameters, we can create an algorithm that performs similar to the specialized ISP algorithm in similar runtime. The construction of such an algorithm is more straightforward as we can create the forward model and embed the workflow in an auto-differentiation framework. This allows us to optimize end-to-end in a straight-forward manner.

## 3.4.4 Optimizing total variation reconstruction and neural networks end-to-end

#### Introduction

In the final case study we will focus on denoising CT reconstructions using convolutional neural networks (CNNs) jointly with total variation reconstruction (TV) [176, 184, 240]. Since neural networks are often large, black-box models that are hard to interpret, it can be desirable to let more computation steps be performed by interpretable algorithms and use a smaller network, rather than using simpler algorithms and a larger neural network. In the experiment we show how embedding CT algorithms in an auto-differentiation framework allows for easy interfacing with deep learning algorithms, and the pipeline can be trained end-to-end. This method allows us to design a more interpretable pipeline using variational methods with similar accuracy, while using a smaller neural network.

Total variation reconstruction is a commonly used variational method to incorporate prior knowledge that the gradient of the image should be sparse. For TV, the regularization term in equation 3.4 becomes the magnitude of the gradient

$$\mathbf{x}^* = \operatorname*{arg\,min}_{\mathbf{x}} \left\{ \|A\mathbf{x} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{x}\|_{\mathrm{TV}} \right\} = \operatorname*{arg\,min}_{\mathbf{x}} \left\{ \|A\mathbf{x} - \mathbf{y}\|_2^2 + \lambda \|\nabla \mathbf{x}\|_1 \right\}.$$
(3.12)

To minimize the functional we use an implementation based on Chambolle-Pock [224]. The regularization parameter  $\lambda$  controls the trade-off between data fidelity and regularization term. A small  $\lambda$  will result in a reconstruction  $\mathbf{x}^*$  that is close to the raw data, but potentially with high noise. A large  $\lambda$  can lead to less noise and more connected components of equal value. For more information on total variation regularization we refer to Rudin-Osher-Fatemi [203].

#### Experiments

We perform an experiment on simulated 2D foam-like phantoms of size  $256 \times 256$ . These phantoms contain both large and small scale features (see Figure 3.16(a)). We simulate parallel beam projections from an angular range of  $-60^{\circ}$  to  $60^{\circ}$ , creating missing wedge artifacts, and add Poisson noise to the projection data (see Figure 3.16(b) for an example FBP reconstruction). For TV reconstruction (Figures 3.16(c) and (d)), this creates a situation where a small  $\lambda$  keeps smaller features, but creates a high noise reconstruction, while a larger  $\lambda$  removes more noise but creates connected components of voxels which removes smaller features. We created a training dataset of 100 randomly generated target phantoms, and corresponding noisy limited angle projections. In our experiments we in each case have the



Figure 3.16: (a) Foam CT binary phantom, (b) FBP reconstruction of phantom projections with added severe noise, (c) TV reconstruction with  $\lambda = 10^{-8}$ , and (d) TV reconstruction with  $\lambda = 10^{-2}$ .

noisy projections as input data, and use the original phantom as target data. Our implementation of total variation reconstruction uses PyTorch primitives, and  $\lambda$  can therefore be optimized end-to-end with gradient-based approaches in conjunction with deep learning algorithms. In total, we trained four pipelines end-to-end.

- 1. **FBP** + **Small CNN:** An FBP reconstruction is input for a small CNN that denoises the image. The CNN consists of 3 layers of  $3 \times 3$  kernels arranged in  $1 \times 64$ ,  $64 \times 32$ , and  $32 \times 1$  channels (19.393 parameters in total). Only the CNN weights are learned.
- 2. Single TV + Small CNN: A single TV reconstruction is input for a small CNN that denoises the image. The CNN consists of 3 layers of  $3 \times 3$  kernels arranged in  $1 \times 64$ ,  $64 \times 32$ , and  $32 \times 1$  channels (19.393 parameters in total). Both  $\lambda$  and the CNN weights are learned jointly.  $\lambda$  is initialized at  $\lambda = 10^{-3}$ .
- 3. Double TV + Small CNN: Two TV reconstructions with different  $\lambda$  are input for a small CNN that uses both inputs to create a single denoised output image. The CNN consists of 3 layers of  $3 \times 3$  kernels arranged in  $2 \times 64$ ,  $64 \times 32$ , and  $32 \times 1$  channels (19.969 parameters in total). Both  $\lambda_1$ ,  $\lambda_2$ , and the CNN weights are learned jointly. The  $\lambda$ 's are initialized as  $\lambda_1 = 10^{-3}$ , and  $\lambda_2 = 10^{-8}$ .
- 4. FBP + Large CNN: An FBP reconstruction is input for a larger CNN that densoises the image. The CNN consists of 4 layers of 3 × 3 kernels arranged in 1 × 160, 160 × 96, 96 × 64, and 64 × 1 channels (195.873 parameters in total). Only the CNN weights are learned.

A diagram of the Double TV + Small CNN pipeline is shown in Figure 3.17. All 4 experiments are run for an equal number of training iterations, and were performed on a workstation with 64GB RAM, an NVidia RTX 2070 Super GPU, and AMD Ryzen 7 3800X CPU. For validating the results we generated an additional random phantom that was not in the training set. The resulting denoised validation reconstructions are shown in Figure 3.18.



Figure 3.17: Pipeline for supervised denoising of CT data.



Phantom (zoom) FBP+small CNN TV+small CNN 2xTV+small CNN FBP+large CNN

Figure 3.18: (Top): From left to right; Foam phantom, output of FBP reconstruction followed by a small CNN, output of a single total variation reconstruction followed by a CNN, output of two total variation reconstructions followed by a CNN, output of FBP reconstruction followed by a larger CNN. (Bottom): Zoom of top row.

We see that all pipelines can denoise the reconstruction significantly, but struggle to remove artifacts introduced by the missing angular information. The FBP+SmallCNN pipeline seemingly performs worse on visual inspection of the zoomed images. Arguably, both TV pipelines perform better than FBP+LargeCNN, even though the larger CNN had 10 times more parameters. This suggests that the prior knowledge incorporated by total variation makes for an easier denoising problem for the CNN. It is unclear whether the addition of a second total variation operator benefited the denoising quality. However, in validation loss the double total variation operator performs better with  $1.805 \cdot 10^{-2}$  loss versus  $1.864 \cdot 10^{-2}$  for the single TV operator. The validation loss for the FBP+SmallCNN pipeline is  $2.322 \cdot 10^{-2}$ , and for the FBP+LargeCNN pipeline  $2.268 \cdot 10^{-2}$ .



Figure 3.19: (a) Mean squared error loss during training. (b) TV regularization  $\lambda$  values during training.

We plot the training losses and  $\lambda$  values during training in Figure 3.19. An interesting observation is that single TV learned a  $\lambda$  with a value in between  $\lambda_1$  and  $\lambda_2$  of the double TV pipeline. We hypothesize that the single TV pipeline had to compromise  $\lambda$  between leaving small features intact, and denoising the reconstruction.

Running total variation in a training procedure comes at large computational cost; a single training step (100 images) took 99.6, and 198.3 seconds for single and double TV respectively whereas FBP + large CNN took 1.82 seconds. For 3D it could therefore be infeasible to train this pipeline. For single scan tuning of  $\lambda$  a surrogate TV approach can be considered during optimization, such as [118]. Alternatively,  $\lambda$  could be tuned on the central slice for 3D cases.

Overall the experiment shows how embedding classical CT algorithms, such as TV reconstruction, in GPU accelerated auto-differentiation frameworks allows for the easy prototyping of mixed classical and deep learning pipelines. We were able to replace a large CNN with 200 thousand parameters by a stack of two TV operators to achieve better denoising accuracy, and reduce the CNN size to 20 thousand parameters. In addition, the total variation pipelines are more interpretable as the behaviour of TV for small or large  $\lambda$  is well understood. In general, this easy interfacing of deep learning and classical methods enables end-to-end learning of parameters, which opens up new areas of research. Additionally, the resulting pipeline may be more interpretable since parameters of classical algorithms are often linked to physical or mathematical concepts that are better understood.

# 3.5 Discussion

Our four use cases and the corresponding experiments demonstrate that a broad range of CT workflows can be implemented as end-to-end optimized pipelines using auto-differentiation. Depending on the particular use case, this approach yields several benefits, which we will now discuss in more detail. When all pieces of a data processing pipeline facilitate gradient propagation, parameters at all steps of the pipeline can be optimized jointly for a criterion calculated at any given step. In combination with explicit quality criteria this allowed us to design workflows where the learnable parameters were used in the earlier stages of the pipeline, while the objective function was more naturally defined at the end of the pipeline. For example, in both the rotation axis alignment and beam hardening correction experiments we were able to optimize parameters that are used in the projection domain for metrics computed in the volume domain. Another benefit is that this approach allows for efficient automatic optimization of parameters that may otherwise be chosen manually. This additionally improves the transferability of CT workflows when applied to new data as the parameters can be optimized in an objective manner using the same quality criterion.

Our approach made it possible to implement workflows more efficiently by using existing building blocks, and by defining the workflow in a declarative manner, i.e., implementing the forward model and then optimizing its parameters with a generic off-the-shelf optimizer. Creating workflows in this manner is typically less time consuming to develop and more flexible compared to specialized methods. For example, in the beam hardening experiment we showed that a gradient descent on the forward model of the physical effect results in a quality comparable to a specialized correction method.

Using our design principles allowed for seamless compatibility between classical and deep learning-based approaches. As classical approaches often come with an intuition of their behaviour, this combination of classical algorithms and deep learning can lead to more robust and interpretable workflows. In the last experiment we showed that using existing classical algorithms in conjunction with deep learning can create workflows that perform similarly to purely deep learning based approaches while using smaller neural networks.

Even though gradient-based end-to-end optimization has several potential benefits, our approach can create certain practical challenges. First, convergence of a gradient descent optimizer is not always guaranteed, and can potentially produce suboptimal solutions. Second, the learning rate needs to be picked manually, which is especially challenging when the involved parameters have significantly different magnitudes. Despite the disadvantages, as we have seen in the field of deep learning, this approach often still produces strong results in practice, and well-performing heuristic solutions to mathematically complicated problems were found.

# 3.6 Conclusion

We have shown how implementing classical CT algorithms in an auto-differentiation framework can improve their transferability and interpretability, enable efficient automatic end-to-end optimization, and allow for solving real-world problems without having to develop specialized algorithms. We have explored four use cases experimentally: rotation axis alignment, phase contrast imaging, beam hardening correction, and end-to-end denoising with deep learning and total variation reconstruction, demonstrating that a wide range of CT workflows can be implemented in such a framework. In the future, the key benefits demonstrated in this paper can be utilized in computational toolboxes that leverage auto-differentiation for improved construction and execution of advanced CT workflows.