



GENERATIONS / VANCOUVER  
12-18 AUGUST  
SIGGRAPH2018

# OPT: A DOMAIN SPECIFIC LANGUAGE FOR NON-LINEAR LEAST SQUARES OPTIMIZATION IN GRAPHICS AND IMAGING

**Zachary DeVito**<sup>1</sup>, **Michael Mara**<sup>1</sup>, Michael Zollhöfer<sup>2</sup>, Gilbert  
Bernstein<sup>1</sup>, Jonathan Ragan-Kelley<sup>3</sup>, Christian Theobalt<sup>2</sup>,  
Pat Hanrahan<sup>1</sup>, Matthew Fisher<sup>4</sup>, Matthias Nießner<sup>5</sup>

1. Stanford University
2. Max-Planck-Institute for Informatics
3. UC Berkeley
4. Adobe Research
5. Technical University of Munich

© 2018 SIGGRAPH. All Rights Reserved



Thank you, lets start right in



## Photography & Recording Encouraged

Recording is allowed for this talk



## LEAST-SQUARES OPTIMIZATION UNDERPINS GRAPHICS

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

Optimization problems are found throughout graphics and vision.

Fundamental techniques like poisson image editing, as-rigid-as-possible warping, and shape from shading are all formulated this way.

At their core, they are just solving least-squares optimization problems over images or meshes.

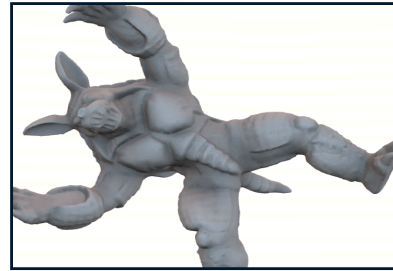
Especially when executed in real-time, these techniques have a bunch of really interesting applications:

## LEAST-SQUARES OPTIMIZATION UNDERPINS GRAPHICS

GENERATIONS / VANCOUVER  
SIGGRAPH 2018



Poisson Image Editing



As-Rigid-As-Possible Deformation



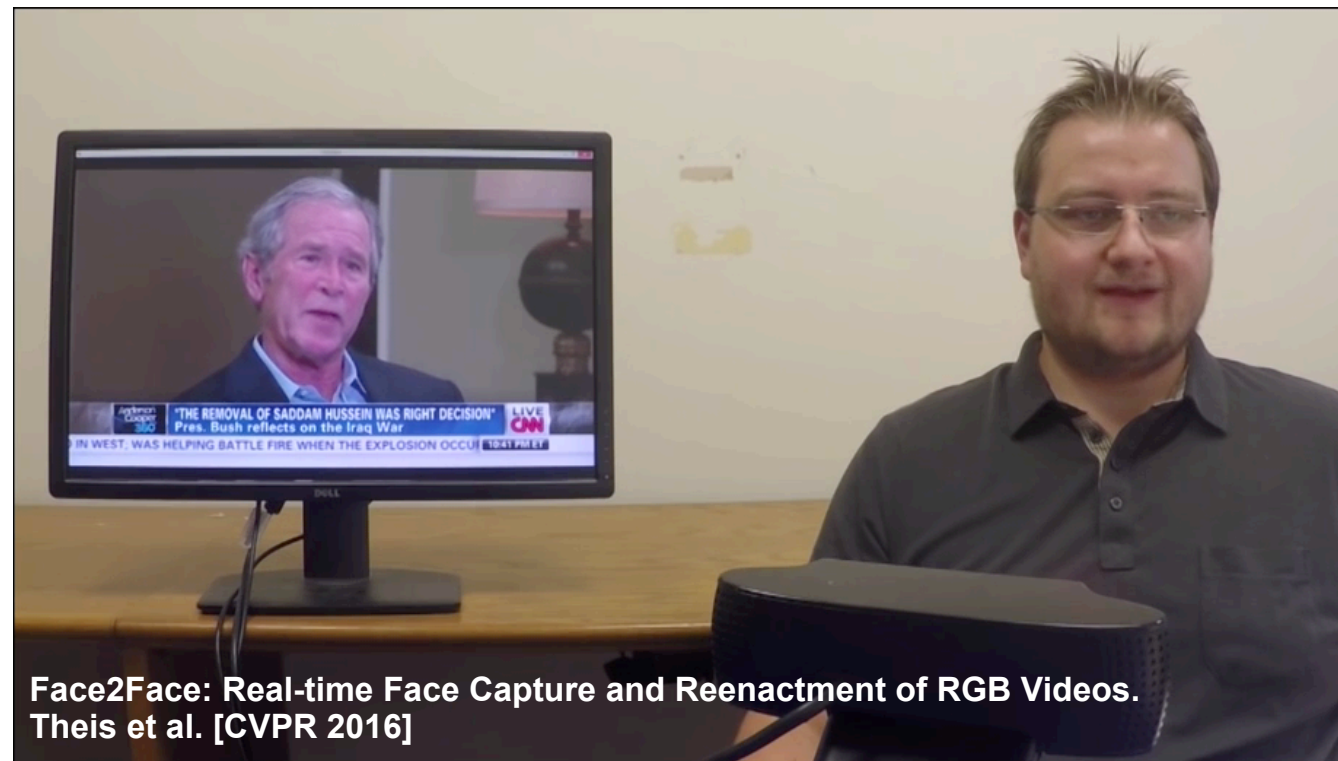
Shape-from-Shading

Optimization problems are found throughout graphics and vision.

Fundamental techniques like poisson image editing, as-rigid-as-possible warping, and shape from shading are all formulated this way.

At their core, they are just solving least-squares optimization problems over images or meshes.

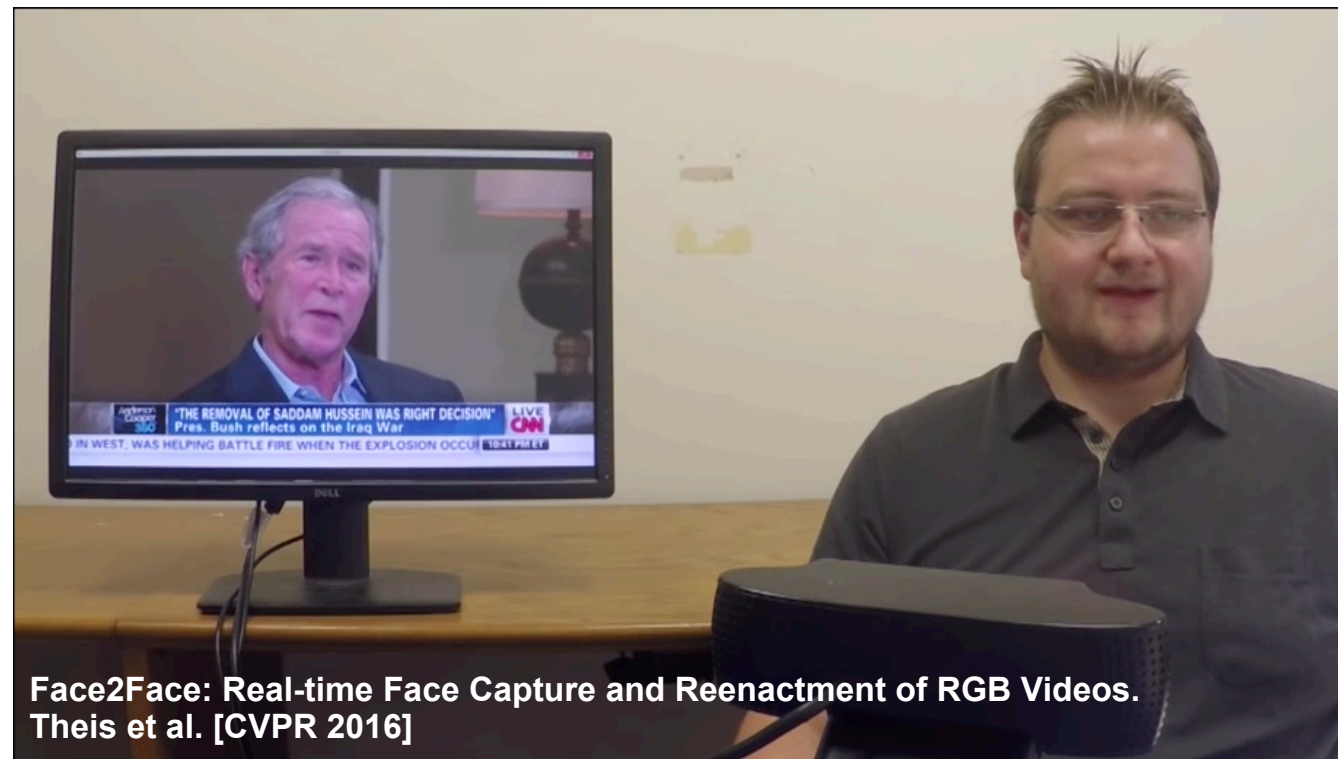
Especially when executed in real-time, these techniques have a bunch of really interesting applications:



**Face2Face: Real-time Face Capture and Reenactment of RGB Videos.**  
Theis et al. [CVPR 2016]

<> You can use webcam to control the facial expressions of a person in a video stream in real-time.

<https://www.youtube.com/watch?v=ohmajJTcpNk>



**Face2Face: Real-time Face Capture and Reenactment of RGB Videos.**  
Theis et al. [CVPR 2016]

<> You can use webcam to control the facial expressions of a person in a video stream in real-time.

<https://www.youtube.com/watch?v=ohmajJTcpNk>

# Relighting



Input RGB



Relit Output



Geometry



Reflectance



Shading



Modified Shading

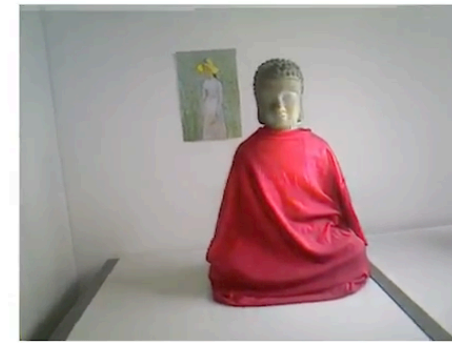
Live User-Guided Intrinsic Video for Static Scenes. Meka et al. [IEEE TVCG 2017]

<> You can decompose scenes into geometry and reflectance, and interactively relight the scene.

# Relighting



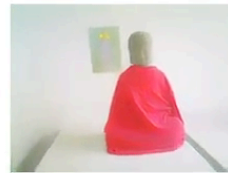
Input RGB



Relit Output



Geometry



Reflectance



Shading



Modified Shading

Live User-Guided Intrinsic Video for Static Scenes. Meka et al. [IEEE TVCG 2017]

<> You can decompose scenes into geometry and reflectance, and interactively relight the scene.

## Material Editing – Plaster to Metal



Input RGB



Output



Reflectance



Shading

Live User-Guided Intrinsic Video for Static Scenes. Meka et al. [IEEE TVCG 2017]

Or you can change the materials of an object in a live scene.

## Material Editing – Plaster to Metal



Input RGB



Output



Reflectance



Shading

Live User-Guided Intrinsic Video for Static Scenes. Meka et al. [IEEE TVCG 2017]

Or you can change the materials of an object in a live scene.



## NONLINEAR LEAST SQUARES ENERGY

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

$$E(\mathbf{x}) = \sum_{r=1}^R [f_r(\mathbf{x})]^2$$

Residual functions

Vector of *unknowns*

Reflectance (*r*)

Shading (*s*)

Input (*i*)

© 2018 SIGGRAPH. All Rights Reserved

7

These problems are often described with non-linear least squared energies, which have this formulation:

There is a vector of unknowns,  $\mathbf{X}$ , which might be pixels in an image or vertices in a graph, and the energy is described use a sum of squared terms,  $f_r$  which are arbitrary functions of the unknowns referred to as residuals.

For instance, when doing relighting our unknowns are the reflectance and shading images. The energy is then a residual per pixel that says the product of the reflectance and shading terms should equal the original image.

## HIGH-LEVEL LIBRARIES MAKE GENERATING SOLVERS EASY

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

**Ceres, CVX, OpenOF, ProxImaL**

If you wanted to get a solver running without much effort, you could use off-the-shelf high-level libraries like Ceres, or CVX

- <> These solvers would use techniques like automatic differentiation to construct explicit sparse matrices,
- <> and would use a sparse matrix library which would run library routines to, for example, solve pcg on the problem to complete an iteration.
- <> This is great from a developer or researchers point of view, since its easy to write new energy functions and try them out
- <> Unfortunately they can be orders of magnitude slower than necessary

## HIGH-LEVEL LIBRARIES MAKE GENERATING SOLVERS EASY

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

### Ceres, CVX, OpenOF, ProxImaL

- CPU Autodiff to construct sparse matrix

If you wanted to get a solver running without much effort, you could use off-the-shelf high-level libraries like Ceres, or CVX

- <> These solvers would use techniques like automatic differentiation to construct explicit sparse matrices,
- <> and would use a sparse matrix library which would run library routines to, for example, solve pcg on the problem to complete an iteration.
- <> This is great from a developer or researchers point of view, since its easy to write new energy functions and try them out
- <> Unfortunately they can be orders of magnitude slower than necessary

## HIGH-LEVEL LIBRARIES MAKE GENERATING SOLVERS EASY

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

### Ceres, CVX, OpenOF, ProxImaL

- CPU Autodiff to construct sparse matrix
- Hand off to sparse-matrix library (e.g. cuSparse) to run PCG loop

If you wanted to get a solver running without much effort, you could use off-the-shelf high-level libraries like Ceres, or CVX

- <> These solvers would use techniques like automatic differentiation to construct explicit sparse matrices,
- <> and would use a sparse matrix library which would run library routines to, for example, solve pcg on the problem to complete an iteration.
- <> This is great from a developer or researchers point of view, since its easy to write new energy functions and try them out
- <> Unfortunately they can be orders of magnitude slower than necessary

## HIGH-LEVEL LIBRARIES MAKE GENERATING SOLVERS EASY

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

### Ceres, CVX, OpenOF, ProxImaL

- CPU Autodiff to construct sparse matrix
- Hand off to sparse-matrix library (e.g. cuSparse) to run PCG loop

(+) easy to write

If you wanted to get a solver running without much effort, you could use off-the-shelf high-level libraries like Ceres, or CVX

- <> These solvers would use techniques like automatic differentiation to construct explicit sparse matrices,
- <> and would use a sparse matrix library which would run library routines to, for example, solve pcg on the problem to complete an iteration.
- <> This is great from a developer or researchers point of view, since its easy to write new energy functions and try them out
- <> Unfortunately they can be orders of magnitude slower than necessary

## HIGH-LEVEL LIBRARIES MAKE GENERATING SOLVERS EASY

GENERATIONS / VANCOUVER  
SIGGRAPH2018

### Ceres, CVX, OpenOF, ProxImaL

- CPU Autodiff to construct sparse matrix
- Hand off to sparse-matrix library (e.g. cuSparse) to run PCG loop

(+) easy to write

(-) significantly slower than optimal

If you wanted to get a solver running without much effort, you could use off-the-shelf high-level libraries like Ceres, or CVX

- <> These solvers would use techniques like automatic differentiation to construct explicit sparse matrices,
- <> and would use a sparse matrix library which would run library routines to, for example, solve pcg on the problem to complete an iteration.
- <> This is great from a developer or researchers point of view, since its easy to write new energy functions and try them out
- <> Unfortunately they can be orders of magnitude slower than necessary

## HANDWRITE TO REACH REALTIME

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

### Per-Energy Custom GPU Solver<sup>[1-6]</sup>

[1] Real-time shading-based refinement for consumer depth cameras Wu et al. Transactions on Graphics 2014  
[2] Real-time non-rigid reconstruction using an RGB-D camera Zollhofer et al. Transactions on Graphics 2014  
[3] Shading-based refinement on volumetric signed distance functions Zollhofer et al. Transactions on Graphics 2015  
[4] Real-time expression transfer for facial reenactment. Thies et al. Transactions on Graphics 2015  
[5] Face2Face: Real-time face capture and reenactment of RGB videos. Thies et al. CVPR 2016  
[6] VolumeDeform: Real-time volumetric non-rigid reconstruction. Innmann et al. ECCV 2016

© 2018 SIGGRAPH. All Rights Reserved

9

In contrast, many real-time techniques, such as written by Wu, Zollhopper, Thies, Innmann and others rely instead on hand-written GPU solvers.

<> These solvers exploit the structure in the images or meshes.

<> They work matrix free, re-constructing needed values on the fly during PCG.

<> And because of this, handwritten derivatives are calculated inside the solvers inner loop.

<> This hand-written approach is incredibly fast, but it is also incredibly hard to get right.

## HANDWRITE TO REACH REALTIME

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

### Per-Energy Custom GPU Solver<sup>[1-6]</sup>

- Exploit image/mesh structure

[1] Real-time shading-based refinement for consumer depth cameras Wu et al. Transactions on Graphics 2014  
[2] Real-time non-rigid reconstruction using an RGB-D camera Zollhofer et al. Transactions on Graphics 2014  
[3] Shading-based refinement on volumetric signed distance functions Zollhofer et al. Transactions on Graphics 2015  
[4] Real-time expression transfer for facial reenactment. Thies et al. Transactions on Graphics 2015  
[5] Face2Face: Real-time face capture and reenactment of RGB videos. Thies et al. CVPR 2016  
[6] VolumeDeform: Real-time volumetric non-rigid reconstruction. Innmann et al. ECCV 2016

© 2018 SIGGRAPH. All Rights Reserved

9

In contrast, many real-time techniques, such as written by Wu, Zollhopper, Thies, Innmann and others rely instead on hand-written GPU solvers.

<> These solvers exploit the structure in the images or meshes.

<> They work matrix free, re-constructing needed values on the fly during PCG.

<> And because of this, handwritten derivatives are calculated inside the solvers inner loop.

<> This hand-written approach is incredibly fast, but it is also incredibly hard to get right.



## HANDWRITE TO REACH REALTIME

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

### Per-Energy Custom GPU Solver<sup>[1-6]</sup>

- Exploit image/mesh structure
- Matrix-free -- Never construct or store system matrices

[1] Real-time shading-based refinement for consumer depth cameras. Wu et al. Transactions on Graphics 2014  
[2] Real-time non-rigid reconstruction using an RGB-D camera. Zollhofer et al. Transactions on Graphics 2014  
[3] Shading-based refinement on volumetric signed distance functions. Zollhofer et al. Transactions on Graphics 2015  
[4] Real-time expression transfer for facial reenactment. Thies et al. Transactions on Graphics 2015  
[5] Face2Face: Real-time face capture and reenactment of RGB videos. Thies et al. CVPR 2016  
[6] VolumeDeform: Real-time volumetric non-rigid reconstruction. Innmann et al. ECCV 2016

© 2018 SIGGRAPH. All Rights Reserved

9

In contrast, many real-time techniques, such as written by Wu, Zollhopper, Thies, Innmann and others rely instead on hand-written GPU solvers.

<> These solvers exploit the structure in the images or meshes.

<> They work matrix free, re-constructing needed values on the fly during PCG.

<> And because of this, handwritten derivatives are calculated inside the solvers inner loop.

<> This hand-written approach is incredibly fast, but its is also incredibly hard to get right.

## HANDWRITE TO REACH REALTIME

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

### Per-Energy Custom GPU Solver<sup>[1-6]</sup>

- Exploit image/mesh structure
- Matrix-free -- Never construct or store system matrices
- Handwritten derivatives are inside the solver kernel

[1] Real-time shading-based refinement for consumer depth cameras Wu et al. Transactions on Graphics 2014  
[2] Real-time non-rigid reconstruction using an RGB-D camera Zollhofer et al. Transactions on Graphics 2014  
[3] Shading-based refinement on volumetric signed distance functions Zollhofer et al. Transactions on Graphics 2015  
[4] Real-time expression transfer for facial reenactment. Thies et al. Transactions on Graphics 2015  
[5] Face2Face: Real-time face capture and reenactment of RGB videos. Thies et al. CVPR 2016  
[6] VolumeDeform: Real-time volumetric non-rigid reconstruction. Innmann et al. ECCV 2016

© 2018 SIGGRAPH. All Rights Reserved

9

In contrast, many real-time techniques, such as written by Wu, Zollhopper, Thies, Innman and others rely instead on hand-written GPU solvers.

<> These solvers exploit the structure in the images or meshes.

<> They work matrix free, re-constructing needed values on the fly during PCG.

<> And because of this, handwritten derivatives are calculated inside the solvers inner loop.

<> This hand-written approach is incredibly fast, but its is also incredibly hard to get right.

## HANDWRITE TO REACH REALTIME

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

### Per-Energy Custom GPU Solver<sup>[1-6]</sup>

- Exploit image/mesh structure
- Matrix-free -- Never construct or store system matrices
- Handwritten derivatives are inside the solver kernel

(+) significantly faster, by orders of magnitude

(-) incredibly hard to write correctly

[1] Real-time shading-based refinement for consumer depth cameras  
Wu et al. Transactions on Graphics 2014

[2] Real-time non-rigid reconstruction using an RGB-D camera  
Zollhofer et al. Transactions on Graphics 2014

[3] Shading-based refinement on volumetric signed distance functions  
Zollhofer et al. Transactions on Graphics 2015

[4] Real-time expression transfer for facial reenactment.  
Thies et al. Transactions on Graphics 2015

[5] Face2Face: Real-time face capture and reenactment of RGB videos. Thies et al. CVPR 2016

[6] VolumeDeform: Real-time volumetric non-rigid reconstruction. Innmann et al. ECCV 2016

© 2018 SIGGRAPH. All Rights Reserved

9

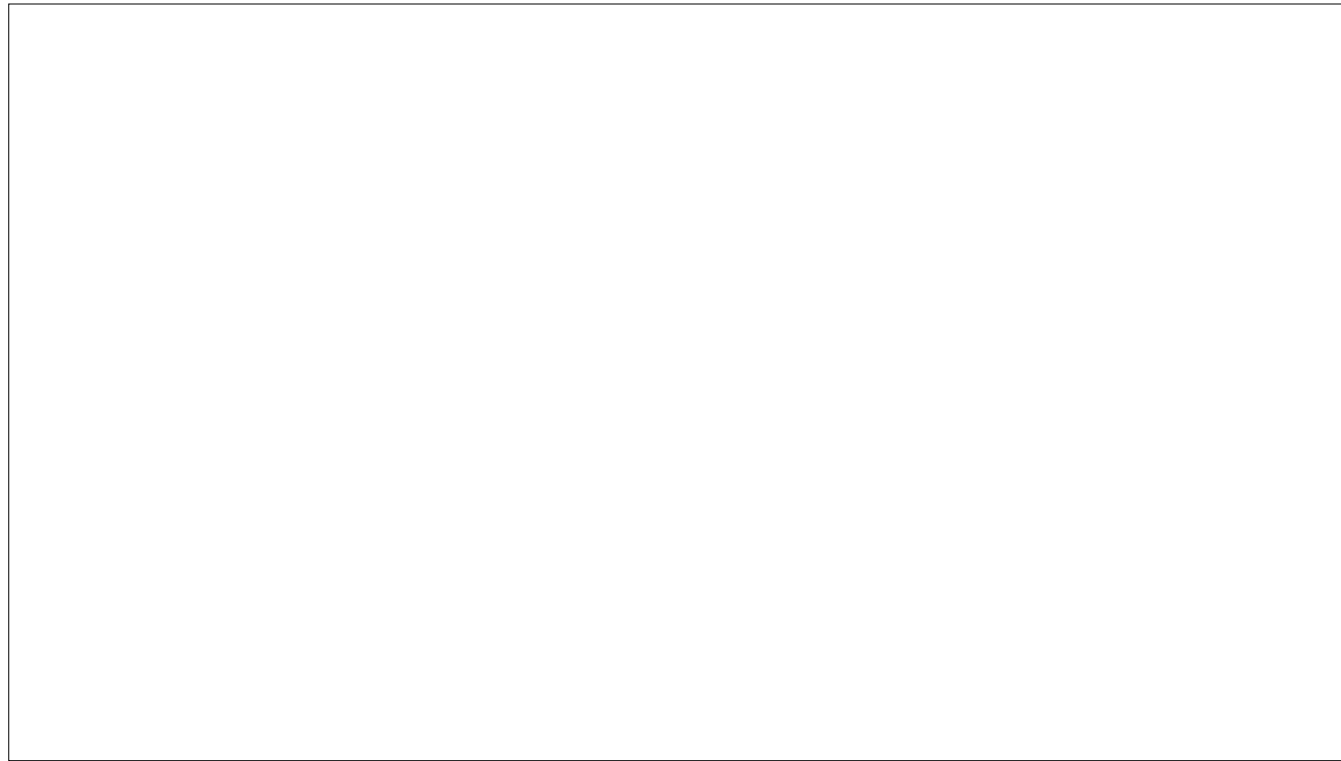
In contrast, many real-time techniques, such as written by Wu, Zollhopper, Thies, Innmann and others rely instead on hand-written GPU solvers.

<> These solvers exploit the structure in the images or meshes.

<> They work matrix free, re-constructing needed values on the fly during PCG.

<> And because of this, handwritten derivatives are calculated inside the solvers inner loop.

<> This hand-written approach is incredibly fast, but it is also incredibly hard to get right.



<> This is code for a real problem. The CUDA code to calculate the energy is pretty concise and completely represents the problem.

<> But to get high-performance, we need all of this code on the right to calculate these in-place matrix products. Writing it by hand is hard, it requires calculus and getting boundary conditions right. The solver code and the energy code is woven together in a complicated way. We never got it right on the first try and bugs would stay in the code for a really long time.

<> This is code for a real problem. The CUDA code to calculate the energy is pretty concise and completely represents the problem.

<> But to get high-performance, we need all of this code on the right to calculate these in-place matrix products. Writing it by hand is hard, it requires calculus and getting boundary conditions right. The solver code and the energy code is woven together in a complicated way. We never got it right on the first try and bugs would stay in the code for a really long time.

# Just the Energy

# Derived hand-written code needed by solver

<> This is code for a real problem. The CUDA code to calculate the energy is pretty concise and completely represents the problem.

<> But to get high-performance, we need all of this code on the right to calculate these in-place matrix products. Writing it by hand is hard, it requires calculus and getting boundary conditions right. The solver code and the energy code is woven together in a complicated way. We never got it right on the first try and bugs would stay in the code for a really long time.

# Productivity vs. Performance

In these existing approaches you have to trade productivity for performance.

<> We set out to build a system that provides both.

# Productivity **and** Performance

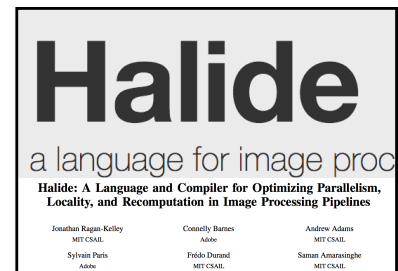
In these existing approaches you have to trade productivity for performance.

<> We set out to build a system that provides both.



# APPROACH: DSLs

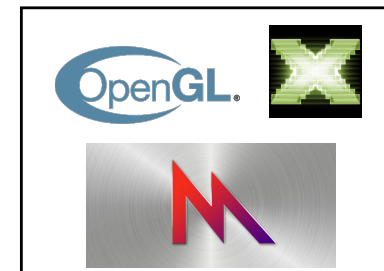
GENERATIONS / VANCOUVER  
SIGGRAPH 2018



Imaging



Simulation

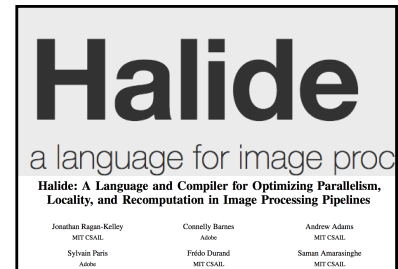


Rendering

In other areas of graphics, such as imaging, physical simulation, or rendering, domain-specific languages have been used to express high-level programs, but still generate high-performance machine code.

# APPROACH: DSLs

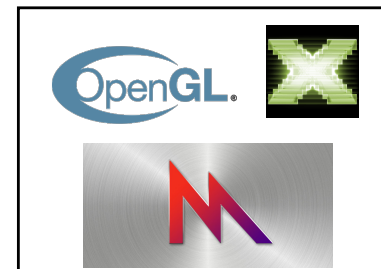
GENERATIONS / VANCOUVER  
SIGGRAPH 2018



Imaging



Simulation



Rendering

In other areas of graphics, such as imaging, physical simulation, or rendering, domain-specific languages have been used to express high-level programs, but still generate high-performance machine code.

## A DSL FOR NLLS OPTIMIZATION

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

Image Warping  
Energy in Opt

```
Offset, Angle = Slice(X,0,2), Slice(X,2,2)
for i,j in Stencil { {1,0}, {-1,0}, {0,1}, {0,-1} } do
  r = (Offset(0,0) - Offset(i,j)) -
    Rotate(Angle(0), OrigPos(0,0) - OrigPos(i,j))
  valid = And(InBounds(i,j),Mask(0,0),Mask(i,j))
  Energy(Select(valid,w_r*r,0))
end
c = w_f*(Offsets(0,0) - Constraints(0,0))
Energy(Select(Valid(Constraints(0,0)),c,0))
```

Opt applies this domain-specific language approach to these optimizations problems. It takes the high-level form of the energy,

<> and automatically produces a real-time GPU solver without all the tedious work.

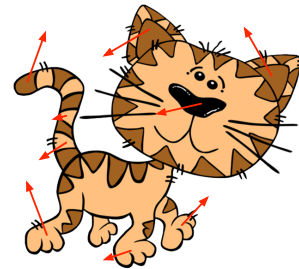
In this talk, we'll first show you what it looks like to write an energy function in Opt, and then walk you through how Opt automatically constructs a gauss-newton style solver from it.

## A DSL FOR NLLS OPTIMIZATION

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

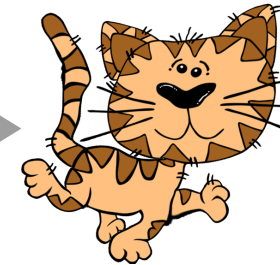
Image Warping  
Energy in Opt

```
Offset, Angle = Slice(X,0,2), Slice(X,2,2)
for i,j in Stencil { {1,0}, {-1,0}, {0,1}, {0,-1} } do
  r = (Offset(0,0) - Offset(i,j)) -
    Rotate(Angle(0), OrigPos(0,0) - OrigPos(i,j))
  valid = And(InBounds(i,j),Mask(0,0),Mask(i,j))
  Energy(Select(valid,w_r*r,0))
end
c = w_f*(Offsets(0,0) - Constraints(0,0))
Energy(Select(Valid(Constraints(0,0)),c,0))
```



Opt Compiler

Fast GPU Solver



Opt applies this domain-specific language approach to these optimizations problems. It takes the high-level form of the energy,

<> and automatically produces a real-time GPU solver without all the tedious work.

In this talk, we'll first show you what it looks like to write an energy function in Opt, and then walk you through how Opt automatically constructs a gauss-newton style solver from it.

## A SIMPLE EXAMPLE

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

$$E(\mathbf{x}) = \sum_{r=1}^R [f_r(\mathbf{x})]^2$$



14

Let's look at a simple Laplacian smoothing problem.

We'll have some notation for our energy term: A is the target input image, and X is the unknown image, which we are trying to find.

<> We can start out with a fitting term, minimizes the difference between X and the original image.

<> Now we can add regularization terms that penalizing the difference between each pixel and its neighbors, blurring the image.

Image from <https://www.flickr.com/photos/cogdog/39525949350/> (public domain)

## A SIMPLE EXAMPLE

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

$$E(\mathbf{x}) = \sum_{r=1}^R [f_r(\mathbf{x})]^2 = \sum_{i=0}^H \sum_{j=0}^W (X_{i,j} - A_{i,j})^2 \quad \text{Fitting Term}$$



Let's look at a simple Laplacian smoothing problem.

We'll have some notation for our energy term: A is the target input image, and X is the unknown image, which we are trying to find.

<> We can start out with a fitting term, minimizes the difference between X and the original image.

<> Now we can add regularization terms that penalizing the difference between each pixel and its neighbors, blurring the image.

Image from <https://www.flickr.com/photos/cogdog/39525949350/> (public domain)

## A SIMPLE EXAMPLE

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

$$E(\mathbf{x}) = \sum_{r=1}^R [f_r(\mathbf{x})]^2 = \sum_{i=0}^H \sum_{j=0}^W (X_{i,j} - A_{i,j})^2 \quad \text{Fitting Term}$$

$$\text{Regularization Terms} \quad + \sum_{i=0}^H \sum_{j=0}^{W-1} (X_{i,j} - X_{i,j+1})^2 + \sum_{i=0}^{H-1} \sum_{j=0}^W (X_{i,j} - X_{i+1,j})^2$$



Let's look at a simple Laplacian smoothing problem.

We'll have some notation for our energy term: A is the target input image, and X is the unknown image, which we are trying to find.

<> We can start out with a fitting term, minimizes the difference between X and the original image.

<> Now we can add regularization terms that penalizing the difference between each pixel and its neighbors, blurring the image.

Image from <https://www.flickr.com/photos/cogdog/39525949350/> (public domain)

# OPT'S PROGRAMMING MODEL

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

```
W,H = Dim("W",0), Dim("H",1)
X = Unknown2D("X",float,{W,H},0)
A = Array2D("A",float,{W,H},1)

w_fit,w_reg = .1,.9
Energy(w_fit*(X(0,0) - A(0,0)), --fitting
       w_reg*(X(0,0) - X(1,0)), --regularization
       w_reg*(X(0,0) - X(0,1)))
```

© 2018 SIGGRAPH. All Rights Reserved

15

Opt allows you to write these energies directly in your problem domain. That is, instead of representing it as a flat list of unknown and residuals, you express the problem in terms of images and meshes.

Here is an Opt version of the Laplacian energy from the last slide.

<> It defines the problem domain it is working on by creating a binding for the original image A, and also an image for the unknown X.

You can have multiple unknown images, and you can also mix images and meshes.

<> Once we have our problem defined in terms of images and meshes, we then define residual energy terms on elements of the domain: this blue term is the same fitting term as the last slide. Note that each term is implicitly squared.

<> These terms are implicitly defined over the entire image like in this illustration.

<> Energies can use a local neighborhood of data using pixel offsets.



## OPT'S PROGRAMMING MODEL

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

Problem Data,  
including Unknowns

```
W,H = Dim("W",0), Dim("H",1)
X = Unknown2D("X",float,{W,H},0)
A = Array2D("A",float,{W,H},1)

w_fit,w_reg = .1,.9
Energy(w_fit*(X(0,0) - A(0,0)), --fitting
       w_reg*(X(0,0) - X(1,0)), --regularization
       w_reg*(X(0,0) - X(0,1)))
```

© 2018 SIGGRAPH. All Rights Reserved

15

Opt allows you to write these energies directly in your problem domain. That is, instead of representing it as a flat list of unknown and residuals, you express the problem in terms of images and meshes.

Here is an Opt version of the Laplacian energy from the last slide.

<> It defines the problem domain it is working on by creating a binding for the original image A, and also an image for the unknown X.

You can have multiple unknown images, and you can also mix images and meshes.

<> Once we have our problem defined in terms of images and meshes, we then define residual energy terms on elements of the domain: this blue term is the same fitting term as the last slide. Note that each term is implicitly squared.

<> These terms are implicitly defined over the entire image like in this illustration.

<> Energies can use a local neighborhood of data using pixel offsets.

## OPT'S PROGRAMMING MODEL

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

```
W,H = Dim("W",0), Dim("H",1)
X = Unknown2D("X",float,{W,H},0)
A = Array2D("A",float,{W,H},1)

w_fit,w_reg = .1,.9
Energy(w_fit*(X(0,0) - A(0,0)), --fitting
       w_reg*(X(0,0) - X(1,0)), --regularization
       w_reg*(X(0,0) - X(0,1)))
```

Problem Data,  
including Unknowns

Residual terms

© 2018 SIGGRAPH. All Rights Reserved

15

Opt allows you to write these energies directly in your problem domain. That is, instead of representing it as a flat list of unknown and residuals, you express the problem in terms of images and meshes.

Here is an Opt version of the Laplacian energy from the last slide.

<> It defines the problem domain it is working on by creating a binding for the original image A, and also an image for the unknown X.

You can have multiple unknown images, and you can also mix images and meshes.

<> Once we have our problem defined in terms of images and meshes, we then define residual energy terms on elements of the domain: this blue term is the same fitting term as the last slide. Note that each term is implicitly squared.

<> These terms are implicitly defined over the entire image like in this illustration.

<> Energies can use a local neighborhood of data using pixel offsets.

# OPT'S PROGRAMMING MODEL

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

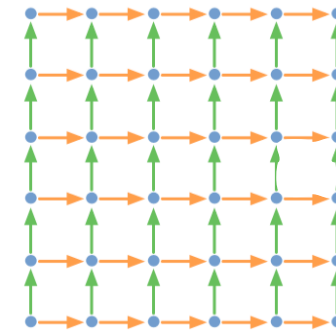
```
W,H = Dim("W",0), Dim("H",1)
X = Unknown2D("X",float,{W,H},0)
A = Array2D("A",float,{W,H},1)
```

Problem Data,  
including Unknowns

```
w_fit,w_reg = .1,.9
Energy(w_fit*(X(0,0) - A(0,0)), --fitting
      w_reg*(X(0,0) - X(1,0)), --regularization
      w_reg*(X(0,0) - X(0,1)))
```

Residual terms

Implied image-wide energies:



© 2018 SIGGRAPH. All Rights Reserved

15

Opt allows you to write these energies directly in your problem domain. That is, instead of representing it as a flat list of unknown and residuals, you express the problem in terms of images and meshes.

Here is an Opt version of the Laplacian energy from the last slide.

<> It defines the problem domain it is working on by creating a binding for the original image A, and also an image for the unknown X.

You can have multiple unknown images, and you can also mix images and meshes.

<> Once we have our problem defined in terms of images and meshes, we then define residual energy terms on elements of the domain: this blue term is the same fitting term as the last slide. Note that each term is implicitly squared.

<> These terms are implicitly defined over the entire image like in this illustration.

<> Energies can use a local neighborhood of data using pixel offsets.

# OPT'S PROGRAMMING MODEL

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

```
W,H = Dim("W",0), Dim("H",1)
X = Unknown2D("X",float,{W,H},0)
A = Array2D("A",float,{W,H},1)
```

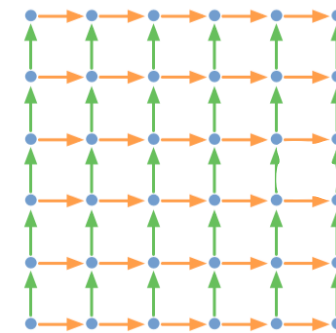
Problem Data,  
including Unknowns

```
w_fit,w_reg = .1,.9
Energy(w_fit*(X(0,0) - A(0,0)), --fitting
      w_reg*(X(0,0) - X(1,0)), --regularization
      w_reg*(X(0,0) - X(0,1)))
```

Residual terms

Pixel Offset

Implied image-wide energies:



© 2018 SIGGRAPH. All Rights Reserved

15

Opt allows you to write these energies directly in your problem domain. That is, instead of representing it as a flat list of unknown and residuals, you express the problem in terms of images and meshes.

Here is an Opt version of the Laplacian energy from the last slide.

<> It defines the problem domain it is working on by creating a binding for the original image A, and also an image for the unknown X.

You can have multiple unknown images, and you can also mix images and meshes.

<> Once we have our problem defined in terms of images and meshes, we then define residual energy terms on elements of the domain: this blue term is the same fitting term as the last slide. Note that each term is implicitly squared.

<> These terms are implicitly defined over the entire image like in this illustration.

<> Energies can use a local neighborhood of data using pixel offsets.

## MINIMIZATION USES DERIVATIVE TERMS

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

$$E(\mathbf{x}) = \sum_{r=1}^R [f_r(\mathbf{x})]^2 = \|\mathbf{F}(\mathbf{x})\|_2^2, \quad \mathbf{F}(\mathbf{x}) = [f_1(\mathbf{x}), \dots, f_R(\mathbf{x})]^T$$

When minimizing one of these energy formulations, you don't compute the direct energy, instead you typically compute terms based on the derivative of the energy.

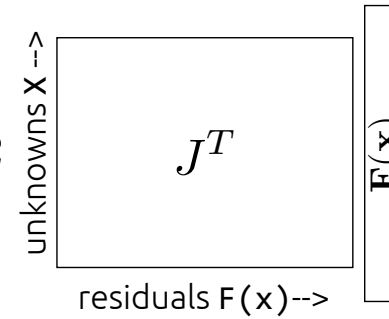
<> In the simplest form, you would compute the gradient, which for least squares problems is this matrix product here,  
<> you can then use gradient descent to step towards the solution, repeatedly calculating this product

## MINIMIZATION USES DERIVATIVE TERMS

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

$$E(\mathbf{x}) = \sum_{r=1}^R [f_r(\mathbf{x})]^2 = \|\mathbf{F}(\mathbf{x})\|_2^2, \quad \mathbf{F}(\mathbf{x}) = [f_1(\mathbf{x}), \dots, f_R(\mathbf{x})]^T$$

$$\nabla E(x) = 2J^T \mathbf{F}(\mathbf{x}) = 2$$



When minimizing one of these energy formulations, you don't compute the direct energy, instead you typically compute terms based on the derivative of the energy.

- <> In the simplest form, you would compute the gradient, which for least squares problems is this matrix product here,
- <> you can then use gradient descent to step towards the solution, repeatedly calculating this product

## MINIMIZATION USES DERIVATIVE TERMS

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

$$E(\mathbf{x}) = \sum_{r=1}^R [f_r(\mathbf{x})]^2 = \|\mathbf{F}(\mathbf{x})\|_2^2, \quad \mathbf{F}(\mathbf{x}) = [f_1(\mathbf{x}), \dots, f_R(\mathbf{x})]^T$$

$$\nabla E(x) = 2J^T \mathbf{F}(\mathbf{x}) = 2 \begin{matrix} \text{unknowns } \mathbf{x} \rightarrow \\ \boxed{J^T} \end{matrix} \begin{matrix} \boxed{\mathbf{F}(\mathbf{x})} \\ \text{residuals } \mathbf{F}(\mathbf{x}) \rightarrow \end{matrix}$$

**Gradient Descent:** Repeat  $\mathbf{x} \leftarrow \mathbf{x} - \alpha 2J^T \mathbf{F}(\mathbf{x})$

When minimizing one of these energy formulations, you don't compute the direct energy, instead you typically compute terms based on the derivative of the energy.

<> In the simplest form, you would compute the gradient, which for least squares problems is this matrix product here,  
<> you can then use gradient descent to step towards the solution, repeatedly calculating this product

## GAUSS-NEWTON IS NEEDED FOR REAL TIME

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

In practice, however, we need higher-order solvers, like Gauss Newton.

These can advance to the solution in fewer steps, allowing us to get to real-time.

In these solvers, an optimization step requires solving a linear system with preconditioned conjugate gradient.

<>The inner-most loop requires calculation of this more complicated matrix product involving the Jacobean matrix  $J$ . Getting this step fast is key to good performance.



## GAUSS-NEWTON IS NEEDED FOR REAL TIME

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

- Gauss-Newton solves a linear system to find  $\Delta \mathbf{x}$  every step
- Use Preconditioned Conjugate Gradient (PCG)

In practice, however, we need higher-order solvers, like Gauss Newton.

These can advance to the solution in fewer steps, allowing us to get to real-time.

In these solvers, an optimization step requires solving a linear system with preconditioned conjugate gradient.

<>The inner-most loop requires calculation of this more complicated matrix product involving the Jacobean matrix  $J$ . Getting this step fast is key to good performance.

## GAUSS-NEWTON IS NEEDED FOR REAL TIME

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

- Gauss-Newton solves a linear system to find  $\Delta \mathbf{x}$  every step
- Use Preconditioned Conjugate Gradient (PCG)

Perf. critical inner loop does repeated calculations of:

$$2J^T J \mathbf{p} = 2 \begin{matrix} \boxed{J^T} & \boxed{J} & \boxed{\mathbf{p}} \end{matrix} \begin{matrix} \leftarrow \text{unknowns } \mathbf{x} \end{matrix}$$

In practice, however, we need higher-order solvers, like Gauss Newton.

These can advance to the solution in fewer steps, allowing us to get to real-time.

In these solvers, an optimization step requires solving a linear system with preconditioned conjugate gradient.

<>The inner-most loop requires calculation of this more complicated matrix product involving the Jacobean matrix J. Getting this step fast is key to good performance.

## OPT IS FAST BECAUSE IT EXPLOITS THE STRUCTURE OF THE PROBLEM DOMAIN

---

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

Typical solvers would explicitly construct the  $J$  matrix, using a sparse matrix format like compressed sparse row, throwing away structure.

The reason Opt can solve so much faster is because it preserves and exploits this structure of the images and meshes that we define the problem on. We can see this by looking at  $J$ , the jacobian matrix

<> The Jacobian matrix  $J$  contain the partial derivatives of each residual with each unknown.

<> For a mesh-based domain you might have unknowns that represent coordinates in space.

## OPT IS FAST BECAUSE IT EXPLOITS THE STRUCTURE OF THE PROBLEM DOMAIN

GENERATIONS / VANCOUVER  
SIGGRAPH 2018


$$J$$

Typical solvers would explicitly construct the  $J$  matrix, using a sparse matrix format like compressed sparse row, throwing away structure.

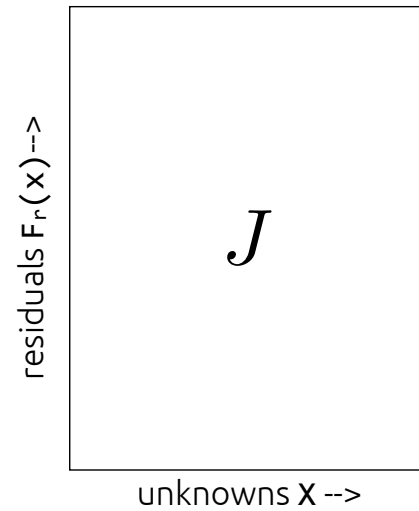
The reason Opt can solve so much faster is because it preserves and exploits this structure of the images and meshes that we define the problem on. We can see this by looking at  $J$ , the jacobian matrix

<> The Jacobian matrix  $J$  contain the partial derivatives of each residual with each unknown.

<> For a mesh-based domain you might have unknowns that represent coordinates in space.

## OPT IS FAST BECAUSE IT EXPLOITS THE STRUCTURE OF THE PROBLEM DOMAIN

GENERATIONS / VANCOUVER  
SIGGRAPH 2018



Typical solvers would explicitly construct the  $J$  matrix, using a sparse matrix format like compressed sparse row, throwing away structure.

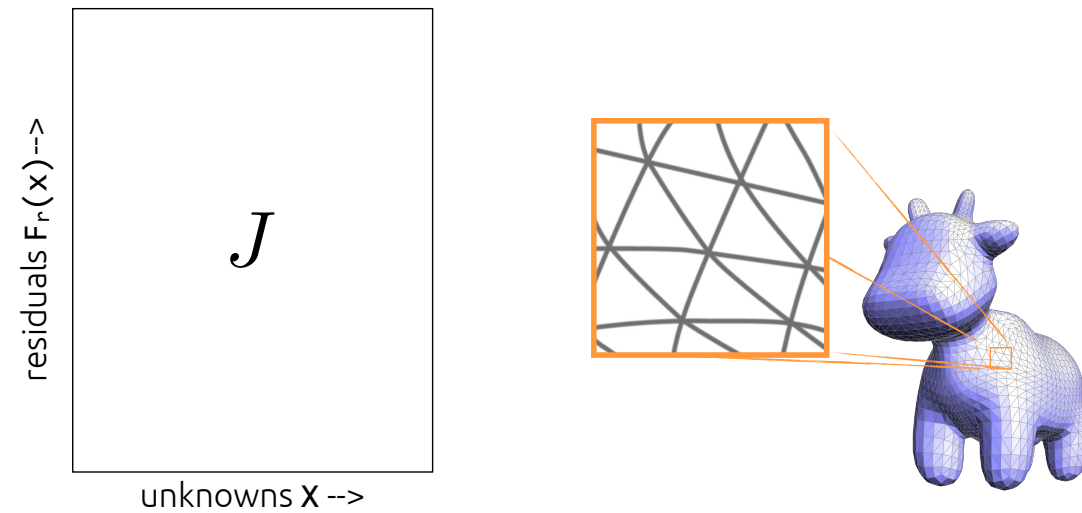
The reason Opt can solve so much faster is because it preserves and exploits this structure of the images and meshes that we define the problem on. We can see this by looking at  $J$ , the jacobian matrix

<> The Jacobian matrix  $J$  contain the partial derivatives of each residual with each unknown.

<> For a mesh-based domain you might have unknowns that represent coordinates in space.

## OPT IS FAST BECAUSE IT EXPLOITS THE STRUCTURE OF THE PROBLEM DOMAIN

GENERATIONS / VANCOUVER  
SIGGRAPH2018



© 2018 SIGGRAPH. All Rights Reserved

18

Typical solvers would explicitly construct the  $J$  matrix, using a sparse matrix format like compressed sparse row, throwing away structure.

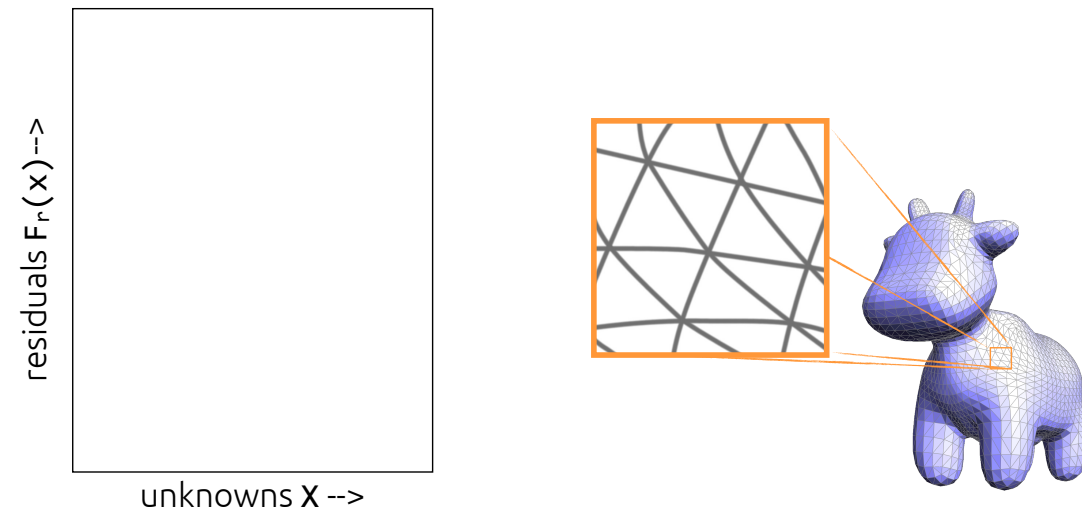
The reason Opt can solve so much faster is because it preserves and exploits this structure of the images and meshes that we define the problem on. We can see this by looking at  $J$ , the jacobian matrix

<> The Jacobian matrix  $J$  contain the partial derivatives of each residual with each unknown.

<> For a mesh-based domain you might have unknowns that represent coordinates in space.

## OPT IS FAST BECAUSE IT EXPLOITS THE STRUCTURE OF THE PROBLEM DOMAIN

GENERATIONS / VANCOUVER  
SIGGRAPH 2018



© 2018 SIGGRAPH. All Rights Reserved

18

Typical solvers would explicitly construct the  $J$  matrix, using a sparse matrix format like compressed sparse row, throwing away structure.

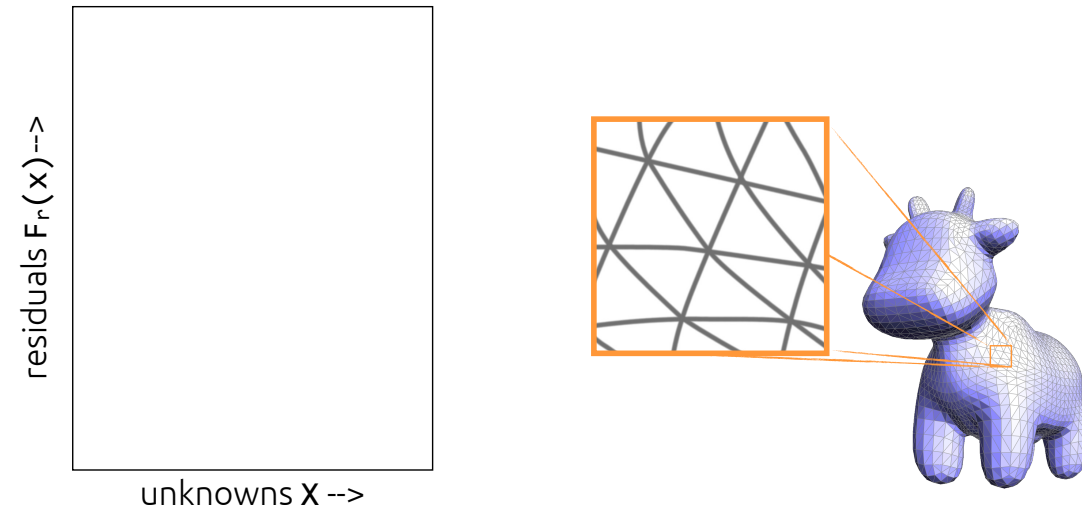
The reason Opt can solve so much faster is because it preserves and exploits this structure of the images and meshes that we define the problem on. We can see this by looking at  $J$ , the jacobian matrix

<> The Jacobian matrix  $J$  contain the partial derivatives of each residual with each unknown.

<> For a mesh-based domain you might have unknowns that represent coordinates in space.

## OPT IS FAST BECAUSE IT EXPLOITS THE STRUCTURE OF THE PROBLEM DOMAIN

GENERATIONS / VANCOUVER  
SIGGRAPH 2018



© 2018 SIGGRAPH. All Rights Reserved

19

This corresponds to three columns of  $J$  per vertex

Then each edge might have several 3D residuals defined,  
<> like a regularization term and a consistency term

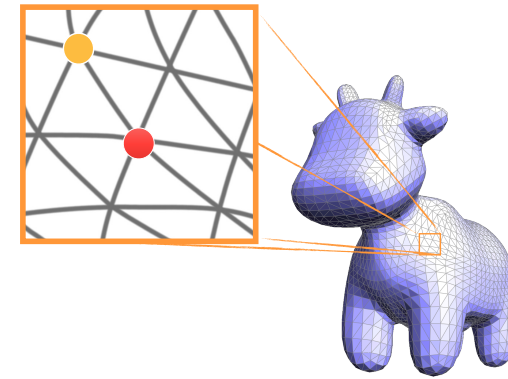
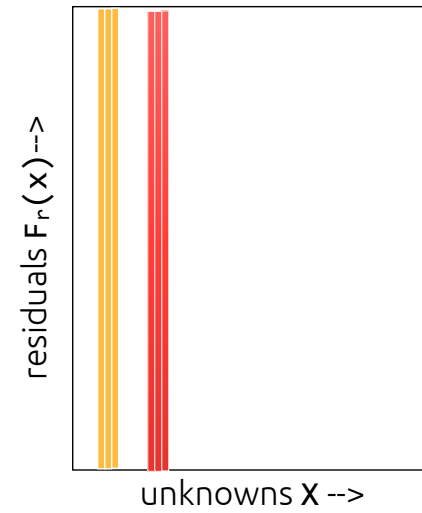
Each of these terms are really 3 residuals and occupy 3 rows of the Jacobian

Because Opt understands the problem domain, it can exploit this structure.



## OPT IS FAST BECAUSE IT EXPLOITS THE STRUCTURE OF THE PROBLEM DOMAIN

GENERATIONS / VANCOUVER  
SIGGRAPH 2018



This corresponds to three columns of  $J$  per vertex

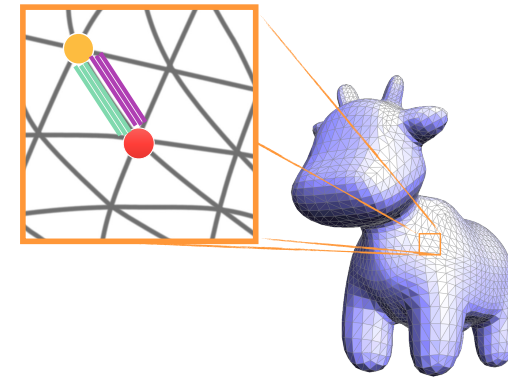
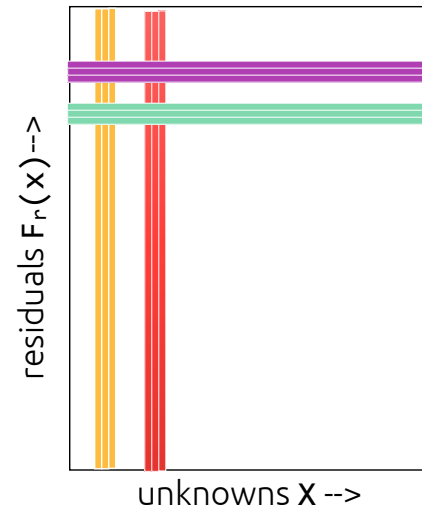
Then each edge might have several 3D residuals defined,  
<> like a regularization term and a consistency term

Each of these terms are really 3 residuals and occupy 3 rows of the Jacobian

Because Opt understands the problem domain, it can exploit this structure.

## OPT IS FAST BECAUSE IT EXPLOITS THE STRUCTURE OF THE PROBLEM DOMAIN

GENERATIONS / VANCOUVER  
SIGGRAPH 2018



This corresponds to three columns of  $J$  per vertex

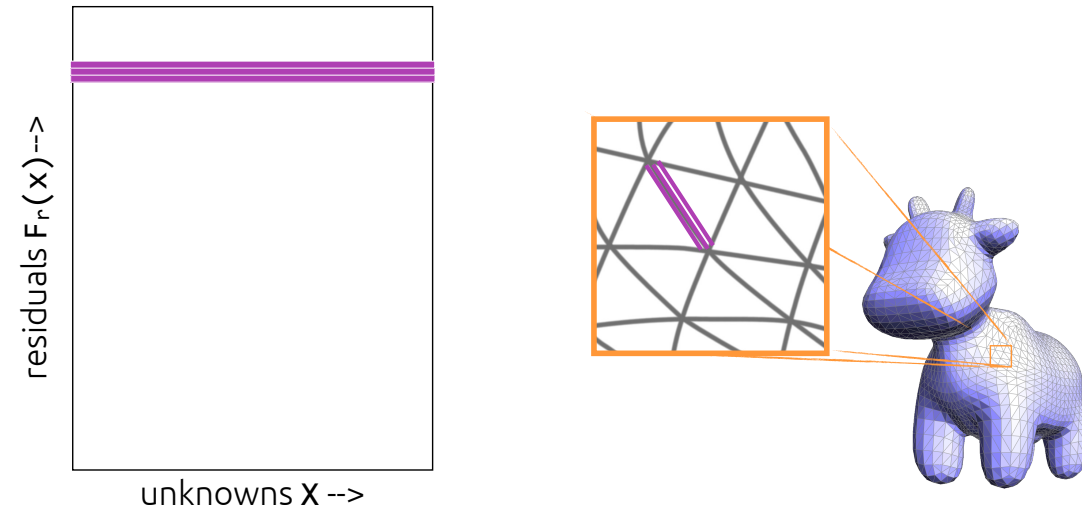
Then each edge might have several 3D residuals defined,  
<> like a regularization term and a consistency term

Each of these terms are really 3 residuals and occupy 3 rows of the Jacobian

Because Opt understands the problem domain, it can exploit this structure.

## OPT IS FAST BECAUSE IT EXPLOITS THE STRUCTURE OF THE PROBLEM DOMAIN

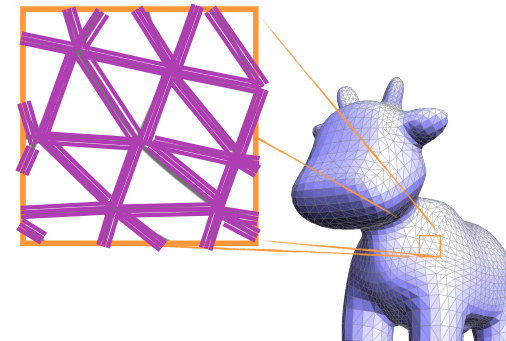
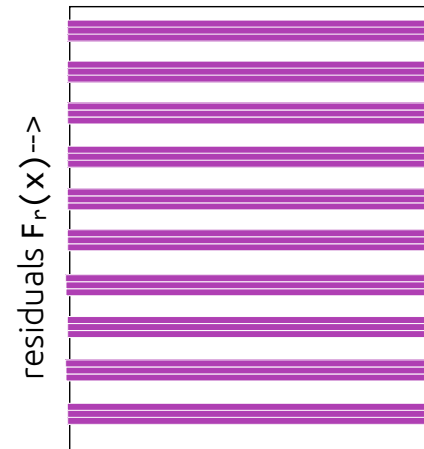
GENERATIONS / VANCOUVER  
SIGGRAPH 2018



First, each residual of a particular type has the same structure, so Opt can calculate terms involving this residual in parallel on the GPU

## OPT IS FAST BECAUSE IT EXPLOITS THE STRUCTURE OF THE PROBLEM DOMAIN

GENERATIONS / VANCOUVER  
SIGGRAPH2018

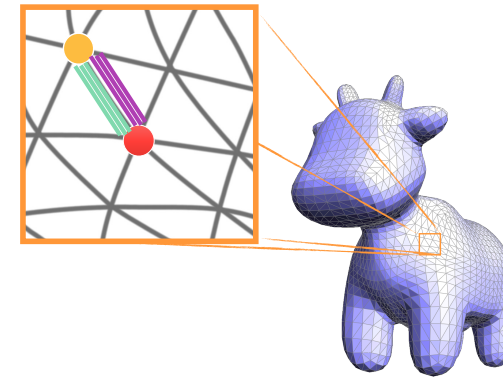
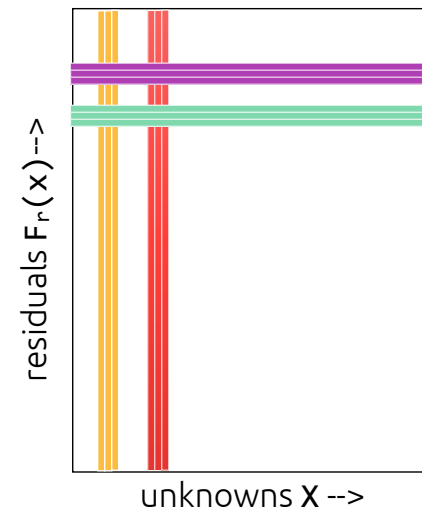


1. **Parallelism:** Each Residual Term Parallelizes by Construction

First, each residual of a particular type has the same structure, so Opt can calculate terms involving this residual in parallel on the GPU

## OPT IS FAST BECAUSE IT EXPLOITS THE STRUCTURE OF THE PROBLEM DOMAIN

GENERATIONS / VANCOUVER  
SIGGRAPH 2018



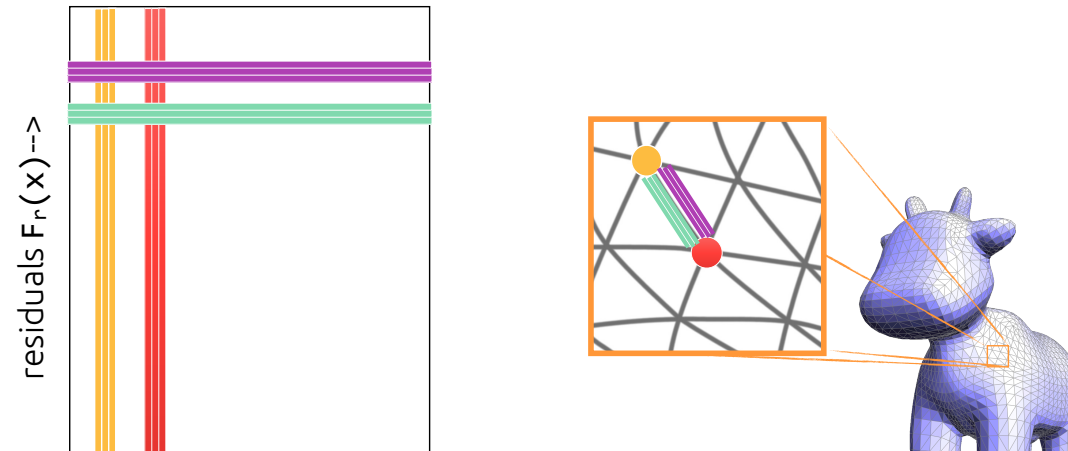
Second, In typical sparse matrices, the connectivity for each row and each column are stored separately.

<> Here for instance, we have six unknowns and six residuals for a total of 36 non-zeros in  $J$

<> But all of these can be derived from a single edge. In the case of images, we don't need to store any connectivity at all, the program defines it!

## OPT IS FAST BECAUSE IT EXPLOITS THE STRUCTURE OF THE PROBLEM DOMAIN

GENERATIONS / VANCOUVER  
SIGGRAPH2018



2. **Connectivity**: One Mesh Edge specifies connectivity for 36 J entries

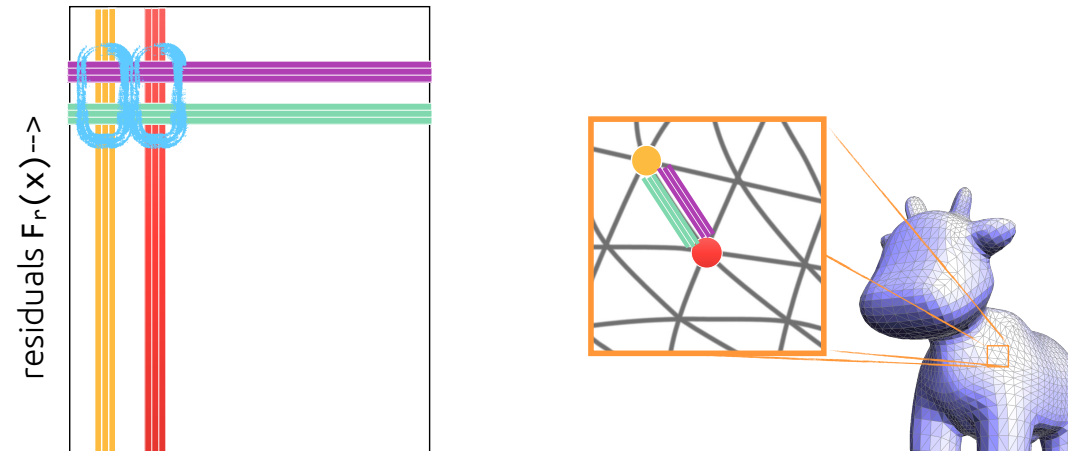
Second, In typical sparse matrices, the connectivity for each row and each column are stored separately.

<> Here for instance, we have six unknowns and six residuals for a total of 36 non-zeros in J

<> But all of these can be derived from a single edge. In the case of images, we don't need to store any connectivity at all, the program defines it!

## OPT IS FAST BECAUSE IT EXPLOITS THE STRUCTURE OF THE PROBLEM DOMAIN

GENERATIONS / VANCOUVER  
SIGGRAPH2018



2. **Connectivity**: One Mesh Edge specifies connectivity for 36 J entries

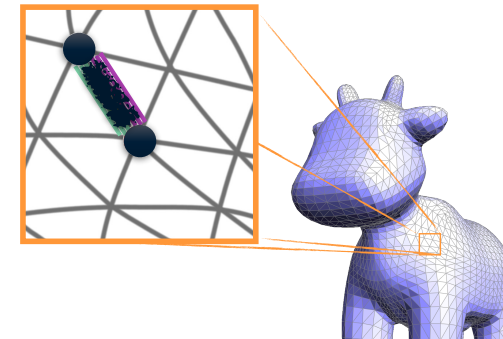
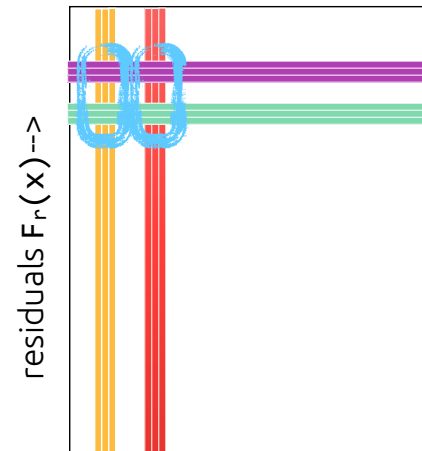
Second, In typical sparse matrices, the connectivity for each row and each column are stored separately.

<> Here for instance, we have six unknowns and six residuals for a total of 36 non-zeros in J

<> But all of these can be derived from a single edge. In the case of images, we don't need to store any connectivity at all, the program defines it!

## OPT IS FAST BECAUSE IT EXPLOITS THE STRUCTURE OF THE PROBLEM DOMAIN

GENERATIONS / VANCOUVER  
SIGGRAPH2018



2. **Connectivity**: One Mesh Edge specifies connectivity for 36 J entries

Second, In typical sparse matrices, the connectivity for each row and each column are stored separately.

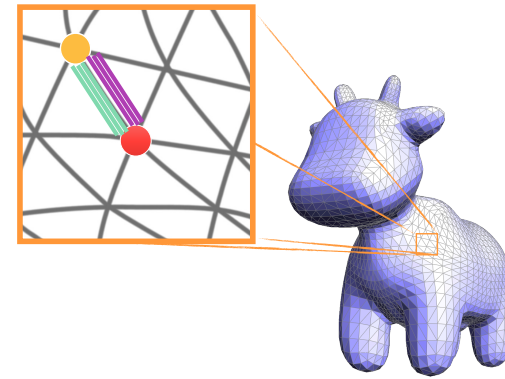
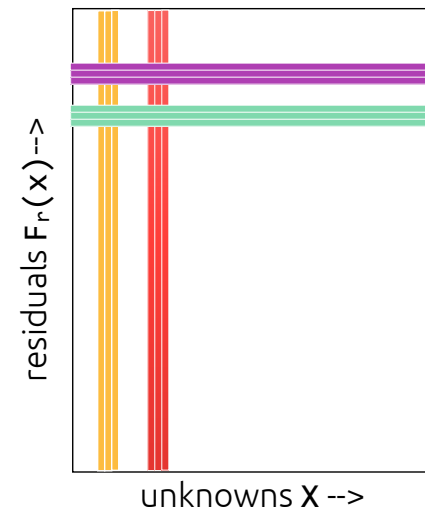
<> Here for instance, we have six unknowns and six residuals for a total of 36 non-zeros in J

<> But all of these can be derived from a single edge. In the case of images, we don't need to store any connectivity at all, the program defines it!



## OPT IS FAST BECAUSE IT EXPLOITS THE STRUCTURE OF THE PROBLEM DOMAIN

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

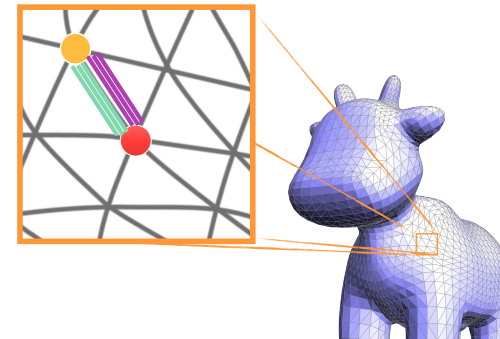
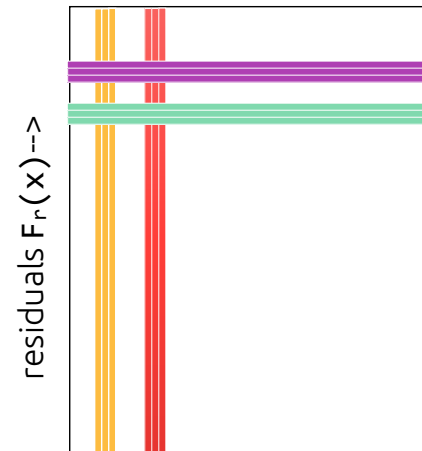


Finally, the 36 non-zeros in  $J$  can be calculated from only 15 problem terms defined on the vertices of the mesh, cutting the amount of memory loaded in half.

By exploiting all three of these properties, we can generate a very fast inner loop.

## OPT IS FAST BECAUSE IT EXPLOITS THE STRUCTURE OF THE PROBLEM DOMAIN

GENERATIONS / VANCOUVER  
SIGGRAPH 2018



3. **Shared Terms:** 2 Vertices' Data (15 floats) turns into 36 J entries

Finally, the 36 non-zeros in  $J$  can be calculated from only 15 problem terms defined on the vertices of the mesh, cutting the amount of memory loaded in half.

By exploiting all three of these properties, we can generate a very fast inner loop.

## Generating the matrix product $g = 2J^T J p$

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

$$\begin{array}{c} \text{↑} \\ \text{unknowns} \end{array} \begin{bmatrix} g_{0,0} \\ \vdots \end{bmatrix} \begin{array}{c} \text{← required row →} \\ \\ \end{array} = 2 \begin{array}{c} \text{↑} \\ \text{unknowns} \\ \text{residuals →} \end{array} \begin{bmatrix} \vdots \\ \vdots \end{bmatrix} \begin{array}{c} \text{↑} \\ \text{residuals} \\ \text{unknowns →} \end{array} \begin{bmatrix} \vdots \\ \vdots \end{bmatrix} \begin{array}{c} \text{↑} \\ \text{unknowns} \end{array} \begin{bmatrix} p \\ \vdots \end{bmatrix}$$

So the challenge in Opt is to generate this matrix product, needed for preconditioned conjugate gradient during Gauss–Newton optimization, from this high-level energy.

\*\*\*Handoff here\*\*\*

<> The structure of this matrix product is specific to a particular set of energy functions.

so we need to use compiler techniques to answer questions using the energy:

<> First, given the energy, we need to find the expression that calculates any particular entry in the Jacobian matrix J.

<> and since J is sparse, we need to use the energy to identify what entries are non-zero.

Both of these questions can be answered using simple program analysis of the energy.

<> For the first problem, we can use a differentiation method to turn the residuals into their partial derivatives.

<> For the second problem, we can use data-dependency analysis that inverts the mapping from residuals to unknowns

## Generating the matrix product $g = 2J^T J p$

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

$$\begin{array}{c} \text{↑} \\ \text{unknowns} \end{array} \begin{bmatrix} g_{0,0} \\ \vdots \end{bmatrix} \begin{array}{c} \text{← required row →} \\ \text{← residuals →} \end{array} = 2 \begin{array}{c} \text{↑} \\ \text{unknowns} \end{array} \begin{bmatrix} \vdots \\ \vdots \end{bmatrix} \begin{array}{c} \text{↑} \\ \text{residuals} \end{array} \begin{bmatrix} \vdots \\ \vdots \end{bmatrix} \begin{array}{c} \text{↑} \\ \text{unknowns} \end{array} \begin{bmatrix} \vdots \\ \vdots \end{bmatrix} \begin{array}{c} \text{↑} \\ \text{unknowns} \end{array} \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}$$

Need compiler techniques to answer two questions using the energy:

So the challenge in Opt is to generate this matrix product, needed for preconditioned conjugate gradient during Gauss–Newton optimization, from this high-level energy.

\*\*\*Handoff here\*\*\*

<> The structure of this matrix product is specific to a particular set of energy functions.

so we need to use compiler techniques to answer questions using the energy:

<> First, given the energy, we need to find the expression that calculates any particular entry in the Jacobian matrix J.

<> and since J is sparse, we need to use the energy to identify what entries are non-zero.

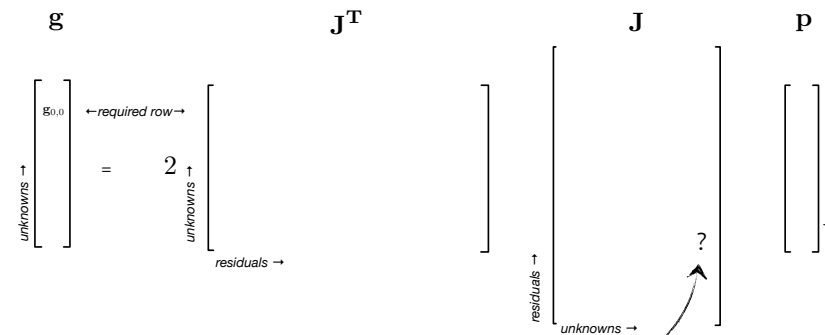
Both of these questions can be answered using simple program analysis of the energy.

<> For the first problem, we can use a differentiation method to turn the residuals into their partial derivatives.

<> For the second problem, we can use data-dependency analysis that inverts the mapping from residuals to unknowns

## Generating the matrix product $g = 2J^T J p$

GENERATIONS / VANCOUVER  
SIGGRAPH 2018



Need compiler techniques to answer two questions using the energy:

1. What is the **expression** to calculate this entry in  $J$ ?

So the challenge in Opt is to generate this matrix product, needed for preconditioned conjugate gradient during Gauss–Newton optimization, from this high-level energy.

\*\*\*Handoff here\*\*\*

<> The structure of this matrix product is specific to a particular set of energy functions.

so we need to use compiler techniques to answer questions using the energy:

<> First, given the energy, we need to find the expression that calculates any particular entry in the Jacobian matrix  $J$ .

<> and since  $J$  is sparse, we need to use the energy to identify what entries are non-zero.

Both of these questions can be answered using simple program analysis of the energy.

<> For the first problem, we can use a differentiation method to turn the residuals into their partial derivatives.

<> For the second problem, we can use data-dependency analysis that inverts the mapping from residuals to unknowns

## Generating the matrix product $g = 2J^T J p$

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

$$\begin{array}{c} \text{↑} \\ \text{unknowns} \end{array} \begin{bmatrix} g_{0,0} \\ \vdots \end{bmatrix} = 2 \begin{array}{c} \text{↑} \\ \text{unknowns} \end{array} \begin{bmatrix} ? & ? & ? & ? \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix} \begin{array}{c} \text{↑} \\ \text{residuals} \end{array} \begin{bmatrix} ? \\ \vdots \end{bmatrix} \begin{array}{c} \text{↑} \\ \text{unknowns} \end{array} \begin{bmatrix} ? \\ \vdots \end{bmatrix}$$

Diagram illustrating the matrix product  $g = 2J^T J p$ . The vector  $g$  (size 1 unknown) is equal to 2 times the product of the transpose of the Jacobian matrix  $J^T$  (size 4 unknowns by 4 residuals) and the product of the Jacobian matrix  $J$  (size 4 residuals by 1 unknown) and the vector  $p$  (size 1 unknown). Arrows indicate the flow of data: from  $J$  to  $J^T$  and from  $J$  to  $p$ .

Need compiler techniques to answer two questions using the energy:

1. What is the **expression** to calculate this entry in  $J$ ?
2. Where are the **non-zero entries** worth examining?

So the challenge in Opt is to generate this matrix product, needed for preconditioned conjugate gradient during Gauss–Newton optimization, from this high-level energy.

\*\*\*Handoff here\*\*\*

<> The structure of this matrix product is specific to a particular set of energy functions.

so we need to use compiler techniques to answer questions using the energy:

<> First, given the energy, we need to find the expression that calculates any particular entry in the Jacobian matrix  $J$ .

<> and since  $J$  is sparse, we need to use the energy to identify what entries are non-zero.

Both of these questions can be answered using simple program analysis of the energy.

<> For the first problem, we can use a differentiation method to turn the residuals into their partial derivatives.

<> For the second problem, we can use data-dependency analysis that inverts the mapping from residuals to unknowns

## Generating the matrix product $g = 2J^T J p$

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

$$\begin{array}{c} \text{↑} \\ \text{unknowns} \end{array} \begin{bmatrix} g_{0,0} \\ \vdots \end{bmatrix} = 2 \begin{array}{c} \text{↑} \\ \text{unknowns} \end{array} \begin{bmatrix} ? & ? & ? & ? \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix} \begin{array}{c} \text{↑} \\ \text{residuals} \end{array} \begin{bmatrix} ? \\ \vdots \end{bmatrix} \begin{array}{c} \text{↑} \\ \text{unknowns} \end{array} \begin{bmatrix} ? \\ \vdots \end{bmatrix}$$

Diagram illustrating the matrix product  $g = 2J^T J p$ . The vector  $g$  (size 1 unknown) is equal to 2 times the product of the Jacobian matrix  $J^T$  (size 4 unknowns by 4 residuals) and the product of the Jacobian matrix  $J$  (size 4 residuals by 1 unknown) and the vector  $p$  (size 1 unknown). Arrows indicate the flow of data: from  $J$  to  $p$ , then to  $J^T$ , and finally to  $g$ . A specific entry in  $J$  is highlighted with a question mark, and an arrow points to it from the text below.

Need compiler techniques to answer two questions using the energy:

1. What is the **expression** to calculate this entry in  $J$ ? *autodiff*
2. Where are the **non-zero entries** worth examining?

So the challenge in Opt is to generate this matrix product, needed for preconditioned conjugate gradient during Gauss–Newton optimization, from this high-level energy.

\*\*\*Handoff here\*\*\*

<> The structure of this matrix product is specific to a particular set of energy functions.

so we need to use compiler techniques to answer questions using the energy:

<> First, given the energy, we need to find the expression that calculates any particular entry in the Jacobian matrix  $J$ .

<> and since  $J$  is sparse, we need to use the energy to identify what entries are non-zero.

Both of these questions can be answered using simple program analysis of the energy.

<> For the first problem, we can use a differentiation method to turn the residuals into their partial derivatives.

<> For the second problem, we can use data-dependency analysis that inverts the mapping from residuals to unknowns

## Generating the matrix product $g = 2J^T J p$

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

$$\begin{array}{c} \text{↑} \\ \text{unknowns} \end{array} \begin{bmatrix} g_{0,0} \\ \vdots \end{bmatrix} = 2 \begin{array}{c} \text{↑} \\ \text{unknowns} \end{array} \begin{bmatrix} ? & ? & ? & ? \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix} \begin{array}{c} \text{↑} \\ \text{residuals} \end{array} \begin{bmatrix} ? \\ \vdots \end{bmatrix} \begin{array}{c} \text{↑} \\ \text{unknowns} \end{array} \begin{bmatrix} ? \\ \vdots \end{bmatrix}$$

Diagram illustrating the matrix product  $g = 2J^T J p$ . The vector  $g$  (size 1 unknown) is equal to 2 times the product of the Jacobian matrix  $J^T$  (size 4 unknowns by 4 residuals) and the product of the Jacobian matrix  $J$  (size 4 residuals by 1 unknown) and the vector  $p$  (size 1 unknown). Arrows indicate the flow of data: from  $J$  to  $p$  to get residuals, then from  $J^T$  and residuals to get  $g$ .

Need compiler techniques to answer two questions using the energy:

1. What is the **expression** to calculate this entry in  $J$ ? *autodiff*
2. Where are the **non-zero entries** worth examining? *data-dependency analysis*

So the challenge in Opt is to generate this matrix product, needed for preconditioned conjugate gradient during Gauss–Newton optimization, from this high-level energy.

\*\*\*Handoff here\*\*\*

<> The structure of this matrix product is specific to a particular set of energy functions.

so we need to use compiler techniques to answer questions using the energy:

<> First, given the energy, we need to find the expression that calculates any particular entry in the Jacobian matrix  $J$ .

<> and since  $J$  is sparse, we need to use the energy to identify what entries are non-zero.

Both of these questions can be answered using simple program analysis of the energy.

<> For the first problem, we can use a differentiation method to turn the residuals into their partial derivatives.

<> For the second problem, we can use data-dependency analysis that inverts the mapping from residuals to unknowns



1. What is the **expression** to calculate this entry in J?

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

```
w_fit, w_reg = .1, .9  
w_fit*(X(0,0) - A(0,0))  
w_reg*(X(0,0) - X(1,0))  
w_reg*(X(0,0) - X(0,1))
```

So the first thing we need to answer is: What is the expression to calculate this entry in J?

<> We transform code written in Opt into a DAG of operators, de-duplicating any reused terms.

At the top is data loaded from the unknown, and at the bottom are individual residual terms at a particular node.

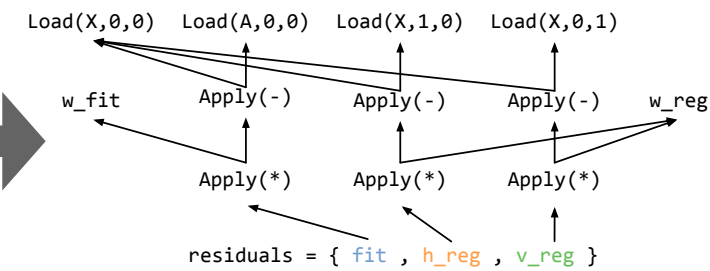
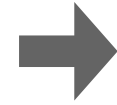
<> Any pair of residual output and unknown input defines one non-zero partial derivative in the J matrix.

<> We can compute the expression for these derivatives using a compiler-based autodiff. In the paper we show how we apply additional optimizations to the resulting expression to minimize the number of terms generated, and handle boundary conditions efficiently.

1. What is the **expression** to calculate this entry in J?

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

```
w_fit, w_reg = .1, .9  
w_fit*(X(0,0) - A(0,0))  
w_reg*(X(0,0) - X(1,0))  
w_reg*(X(0,0) - X(0,1))
```



So the first thing we need to answer is: What is the expression to calculate this entry in J?

<> We transform code written in Opt into a DAG of operators, de-duplicating any reused terms.

At the top is data loaded from the unknown, and at the bottom are individual residual terms at a particular node.

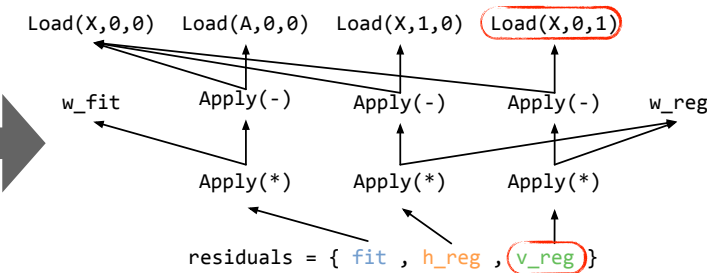
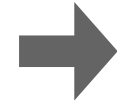
<> Any pair of residual output and unknown input defines one non-zero partial derivative in the J matrix.

<> We can compute the expression for these derivatives using a compiler-based autodiff. In the paper we show how we apply additional optimizations to the resulting expression to minimize the number of terms generated, and handle boundary conditions efficiently.

1. What is the **expression** to calculate this entry in J?

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

```
w_fit, w_reg = .1, .9  
w_fit*(X(0,0) - A(0,0))  
w_reg*(X(0,0) - X(1,0))  
w_reg*(X(0,0) - X(0,1))
```



So the first thing we need to answer is: What is the expression to calculate this entry in J?

<> We transform code written in Opt into a DAG of operators, de-duplicating any reused terms.

At the top is data loaded from the unknown, and at the bottom are individual residual terms at a particular node.

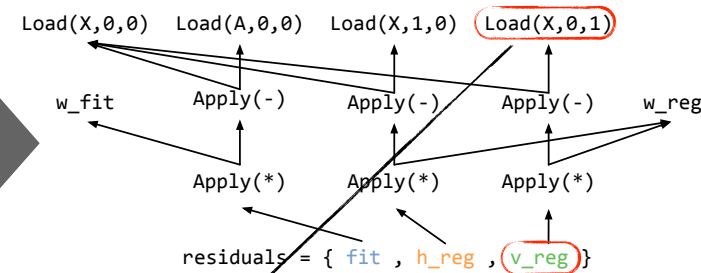
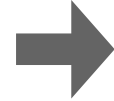
<> Any pair of residual output and unknown input defines one non-zero partial derivative in the J matrix.

<> We can compute the expression for these derivatives using a compiler-based autodiff. In the paper we show how we apply additional optimizations to the resulting expression to minimize the number of terms generated, and handle boundary conditions efficiently.

1. What is the **expression** to calculate this entry in J?

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

```
w_fit, w_reg = .1, .9  
w_fit*(X(0,0) - A(0,0))  
w_reg*(X(0,0) - X(1,0))  
w_reg*(X(0,0) - X(0,1))
```



```
gen_derivative(v_reg, X(0,1)) --> -w_reg
```

So the first thing we need to answer is: What is the expression to calculate this entry in J?

<> We transform code written in Opt into a DAG of operators, de-duplicating any reused terms.

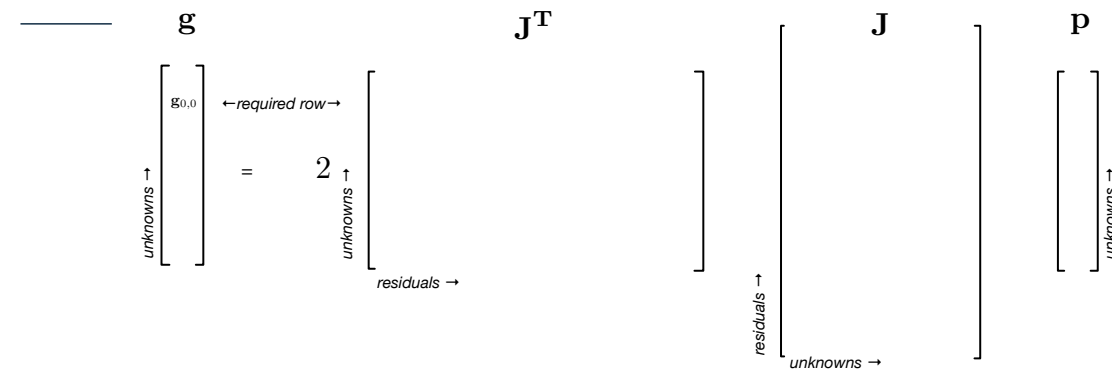
At the top is data loaded from the unknown, and at the bottom are individual residual terms at a particular node.

<> Any pair of residual output and unknown input defines one non-zero partial derivative in the J matrix.

<> We can compute the expression for these derivatives using a compiler-based autodiff. In the paper we show how we apply additional optimizations to the resulting expression to minimize the number of terms generated, and handle boundary conditions efficiently.

## 2. Where are the **non-zero entries** worth examining?

GENERATIONS / VANCOUVER  
SIGGRAPH 2018



Our second problem is to find the non-zeros needed to calculate the matrix products.

<> For a particular matrix product, We work left to right, identifying the non-zeros we need.

For example, in this row of J transpose, we ask

“what are the non-zero columns related this unknown”

<> These are the columns corresponding to the residuals that use that unknown.

<> to do this we need to compute a mapping from an unknown to the residuals that use it.

<> We can identify non-zeros in J as well.

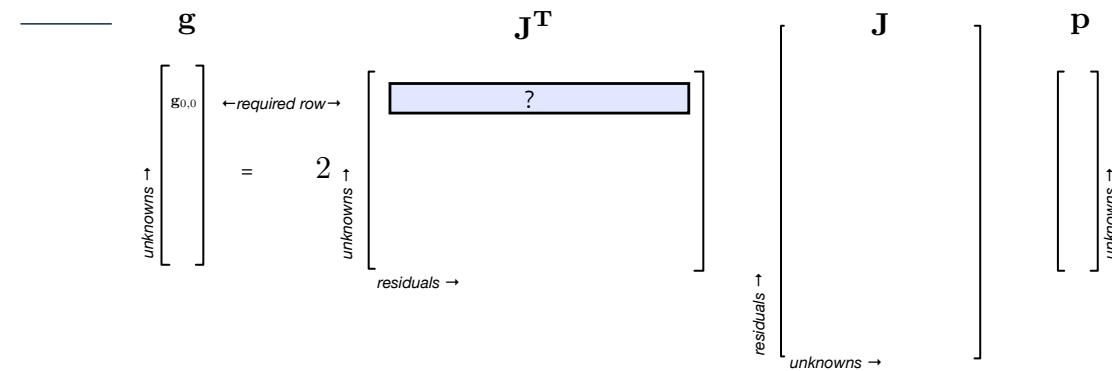
<> This is a very similar problem, but because of the transpose we have one non-zero column for each unknown used by a residual.

<> This requires a map from residual to unknown.

In the image case, we can derive this information from the stencil access patterns, and in the mesh cases, we can recover it from the connectivity information itself. We provide more details in the paper.

## 2. Where are the **non-zero entries** worth examining?

GENERATIONS / VANCOUVER  
SIGGRAPH 2018



For this row, corresponding to  $X_{00}$  what are the non-zero columns?

Our second problem is to find the non-zeros needed to calculate the matrix products.

<> For a particular matrix product, We work left to right, identifying the non-zeros we need.

For example, in this row of  $\mathbf{J}$  transpose, we ask

“what are the non-zero columns related this unknown”

<> These are the columns corresponding to the residuals that use that unknown.

<> to do this we need to compute a mapping from an unknown to the residuals that use it.

<> We can identify non-zeros in  $\mathbf{J}$  as well.

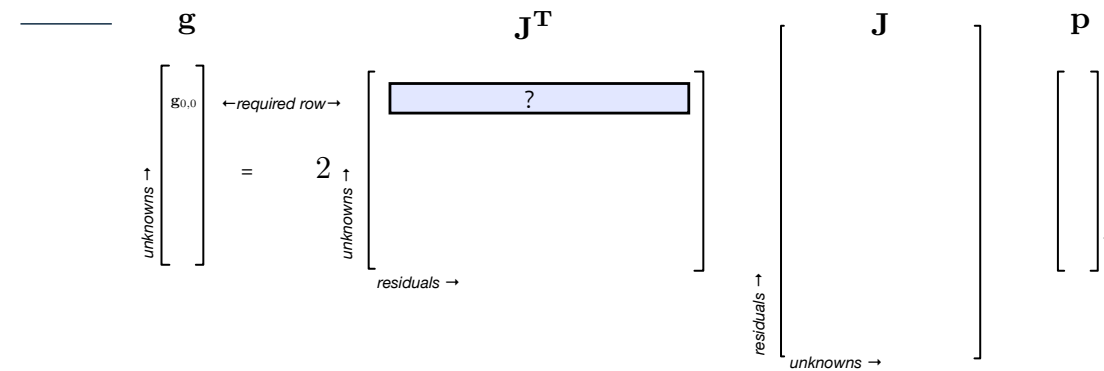
<> This is a very similar problem, but because of the transpose we have one non-zero column for each unknown used by a residual.

<> This requires a map from residual to unknown.

In the image case, we can derive this information from the stencil access patterns, and in the mesh cases, we can recover it from the connectivity information itself. We provide more details in the paper.

## 2. Where are the **non-zero entries** worth examining?

GENERATIONS / VANCOUVER  
SIGGRAPH 2018



For this row, corresponding to  $X_{00}$  what are the non-zero columns?

one for each residual term  $f_r(\mathbf{x})$  that uses  $X_{00}$

Our second problem is to find the non-zeros needed to calculate the matrix products.

<> For a particular matrix product, We work left to right, identifying the non-zeros we need.

For example, in this row of  $\mathbf{J}$  transpose, we ask

“what are the non-zero columns related this unknown”

<> These are the columns corresponding to the residuals that use that unknown.

<> to do this we need to compute a mapping from an unknown to the residuals that use it.

<> We can identify non-zeros in  $\mathbf{J}$  as well.

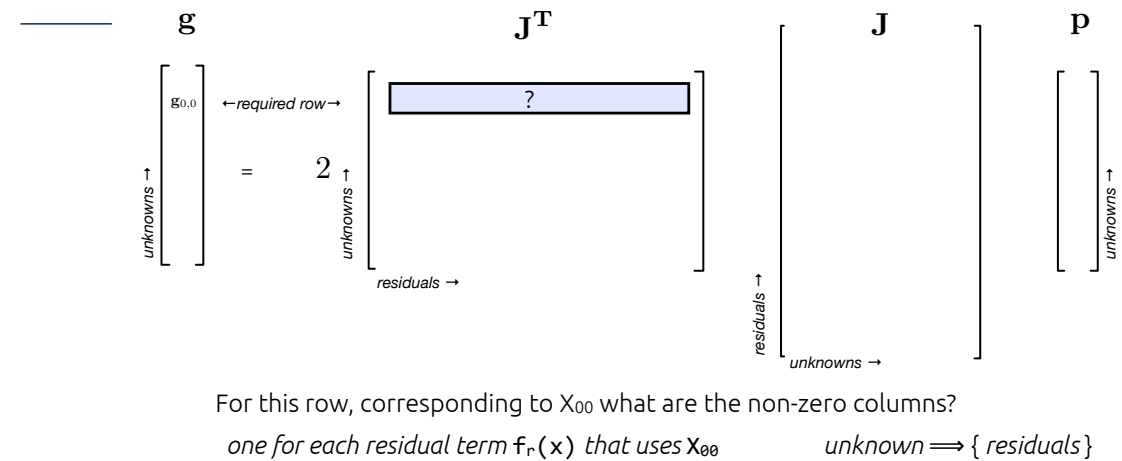
<> This is a very similar problem, but because of the transpose we have one non-zero column for each unknown used by a residual.

<> This requires a map from residual to unknown.

In the image case, we can derive this information from the stencil access patterns, and in the mesh cases, we can recover it from the connectivity information itself. We provide more details in the paper.

## 2. Where are the **non-zero entries** worth examining?

GENERATIONS / VANCOUVER  
SIGGRAPH 2018



© 2018 SIGGRAPH. All Rights Reserved

25

Our second problem is to find the non-zeros needed to calculate the matrix products.

<> For a particular matrix product, We work left to right, identifying the non-zeros we need.

For example, in this row of J transpose, we ask

“what are the non-zero columns related this unknown”

<> These are the columns corresponding to the residuals that use that unknown.

<> to do this we need to compute a mapping from an unknown to the residuals that use it.

<> We can identify non-zeros in J as well.

<> This is a very similar problem, but because of the transpose we have one non-zero column for each unknown used by a residual.

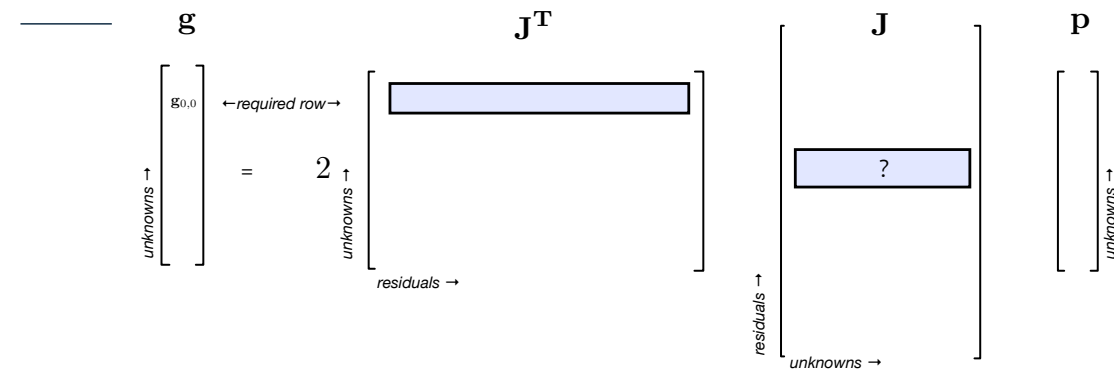
<> This requires a map from residual to unknown.

In the image case, we can derive this information from the stencil access patterns, and in the mesh cases, we can recover it from the connectivity information itself. We provide more details in the paper.



## 2. Where are the **non-zero entries** worth examining?

GENERATIONS / VANCOUVER  
SIGGRAPH 2018



For this row, corresponding to  $X_{00}$  what are the non-zero columns?

one for each residual term  $f_r(x)$  that uses  $X_{00}$        $unknown \Rightarrow \{residuals\}$

For this row, corresponding to  $R_{00}$  what are the non-zero columns?

Our second problem is to find the non-zeros needed to calculate the matrix products.

<> For a particular matrix product, We work left to right, identifying the non-zeros we need.

For example, in this row of  $J$  transpose, we ask

“what are the non-zero columns related this unknown”

<> These are the columns corresponding to the residuals that use that unknown.

<> to do this we need to compute a mapping from an unknown to the residuals that use it.

<> We can identify non-zeros in  $J$  as well.

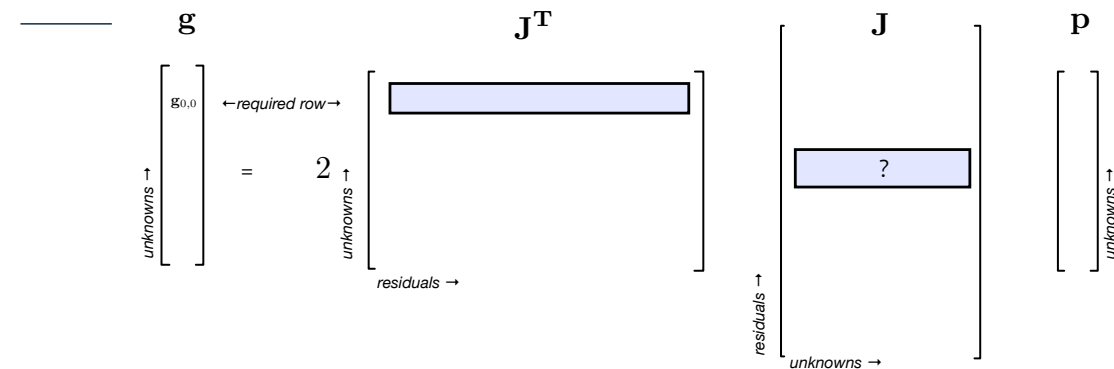
<> This is a very similar problem, but because of the transpose we have one non-zero column for each unknown used by a residual.

<> This requires a map from residual to unknown.

In the image case, we can derive this information from the stencil access patterns, and in the mesh cases, we can recover it from the connectivity information itself. We provide more details in the paper.

## 2. Where are the **non-zero entries** worth examining?

GENERATIONS / VANCOUVER  
SIGGRAPH 2018



For this row, corresponding to  $X_{00}$  what are the non-zero columns?

one for each residual term  $f_r(x)$  that uses  $X_{00}$        $unknown \Rightarrow \{residuals\}$

For this row, corresponding to  $R_{00}$  what are the non-zero columns?

one for each unknown term  $x_{ij}$  used by  $f_{r00}(x)$

© 2018 SIGGRAPH. All Rights Reserved

25

Our second problem is to find the non-zeros needed to calculate the matrix products.

<> For a particular matrix product, We work left to right, identifying the non-zeros we need.

For example, in this row of  $J$  transpose, we ask

“what are the non-zero columns related this unknown”

<> These are the columns corresponding to the residuals that use that unknown.

<> to do this we need to compute a mapping from an unknown to the residuals that use it.

<> We can identify non-zeros in  $J$  as well.

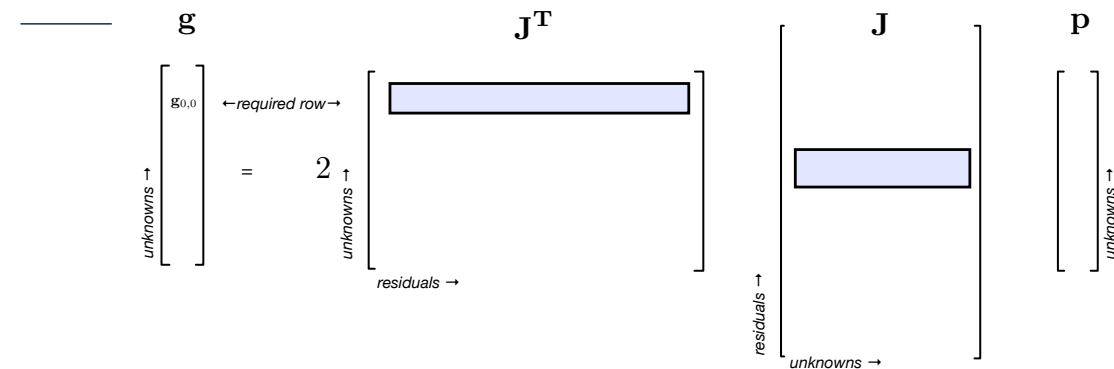
<> This is a very similar problem, but because of the transpose we have one non-zero column for each unknown used by a residual.

<> This requires a map from residual to unknown.

In the image case, we can derive this information from the stencil access patterns, and in the mesh cases, we can recover it from the connectivity information itself. We provide more details in the paper.

## 2. Where are the **non-zero entries** worth examining?

GENERATIONS / VANCOUVER  
SIGGRAPH 2018



For this row, corresponding to  $X_{00}$  what are the non-zero columns?

one for each residual term  $f_r(\mathbf{x})$  that uses  $X_{00}$        $\text{unknown} \Rightarrow \{\text{residuals}\}$

For this row, corresponding to  $R_{00}$  what are the non-zero columns?

one for each unknown term  $x_{ij}$  used by  $f_{r00}(\mathbf{x})$        $\text{residual} \Rightarrow \{\text{unknowns}\}$

© 2018 SIGGRAPH. All Rights Reserved

25

Our second problem is to find the non-zeros needed to calculate the matrix products.

<> For a particular matrix product, We work left to right, identifying the non-zeros we need.

For example, in this row of  $\mathbf{J}$  transpose, we ask

“what are the non-zero columns related this unknown”

<> These are the columns corresponding to the residuals that use that unknown.

<> to do this we need to compute a mapping from an unknown to the residuals that use it.

<> We can identify non-zeros in  $\mathbf{J}$  as well.

<> This is a very similar problem, but because of the transpose we have one non-zero column for each unknown used by a residual.

<> This requires a map from residual to unknown.

In the image case, we can derive this information from the stencil access patterns, and in the mesh cases, we can recover it from the connectivity information itself. We provide more details in the paper.

# PUTTING IT ALL TOGETHER

---

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

Finally we compose the pieces.

<> We use our unknown<->residual mappings to find non-zero partial derivatives and come up with a matrix-free equation for a single entry of the output, treating the derivatives as placeholders

<> We then use compile-time auto diff to generate an expression for each non-zero partial derivative, slotting them into the equation from step 1.

<> Finally, we parallelize across outputs of the matrix product. For large problems this easily saturates even high-end GPUs.

# PUTTING IT ALL TOGETHER

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

1. Where are the non-zeros in this expression?

*Use dependency analysis to find them, and generate multiplication expression using them.*

$$g_{0,0} = 2 \frac{dft_{0,0}}{dx_{0,0}} \frac{dft_{0,0}}{dx_{0,0}} p_{0,0} + 2 \underbrace{\frac{dh_{reg_{0,0}}}{dx_{0,0}}}_{\text{from } \mathbf{J}^T} \left( \underbrace{\frac{dh_{reg_{0,0}}}{dx_{0,0}} p_{0,0} + \frac{dh_{reg_{0,0}}}{dx_{1,0}} p_{1,0}}_{\text{from } \mathbf{J}} \right) + \dots$$

Finally we compose the pieces.

<> We use our unknown<->residual mappings to find non-zero partial derivatives and come up with a matrix-free equation for a single entry of the output, treating the derivatives as placeholders

<> We then use compile-time auto diff to generate an expression for each non-zero partial derivative, slotting them into the equation from step 1.

<> Finally, we parallelize across outputs of the matrix product. For large problems this easily saturates even high-end GPUs.

# PUTTING IT ALL TOGETHER

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

1. Where are the non-zeros in this expression?

*Use dependency analysis to find them, and generate multiplication expression using them.*

$$g_{0,0} = 2 \frac{df_{0,0}}{dx_{0,0}} \frac{df_{0,0}}{dx_{0,0}} p_{0,0} + 2 \underbrace{\frac{dh_{reg_{0,0}}}{dx_{0,0}}}_{\text{from } J^T} \left( \underbrace{\frac{dh_{reg_{0,0}}}{dx_{0,0}}}_{\text{from } J} p_{0,0} + \underbrace{\frac{dh_{reg_{0,0}}}{dx_{1,0}}}_{\text{from } J} p_{1,0} \right) + \dots$$

2. What are the values of the non-zeroes?

*Use automatic differentiation.*

$$= 2w_{fit}^2 p_{0,0} + 2w_{reg}(w_{reg} p_{0,0} + -w_{reg} p_{1,0}) + \dots$$

Finally we compose the pieces.

<> We use our unknown<->residual mappings to find non-zero partial derivatives and come up with a matrix-free equation for a single entry of the output, treating the derivatives as placeholders

<> We then use compile-time auto diff to generate an expression for each non-zero partial derivative, slotting them into the equation from step 1.

<> Finally, we parallelize across outputs of the matrix product. For large problems this easily saturates even high-end GPUs.

# PUTTING IT ALL TOGETHER

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

1. Where are the non-zeros in this expression?

*Use dependency analysis to find them, and generate multiplication expression using them.*

$$g_{0,0} = 2 \frac{df_{0,0}}{dx_{0,0}} \frac{df_{0,0}}{dx_{0,0}} p_{0,0} + 2 \underbrace{\frac{dh_{reg_{0,0}}}{dx_{0,0}}}_{\text{from } J^T} \left( \underbrace{\frac{dh_{reg_{0,0}}}{dx_{0,0}} p_{0,0}}_{\text{from } J} + \underbrace{\frac{dh_{reg_{0,0}}}{dx_{1,0}} p_{1,0}}_{\text{from } J} \right) + \dots$$

2. What are the values of the non-zeroes?

*Use automatic differentiation.*

$$= 2w_{fit}^2 p_{0,0} + 2w_{reg}(w_{reg} p_{0,0} + -w_{reg} p_{1,0}) + \dots$$

3. Parallelize across outputs.

Finally we compose the pieces.

<> We use our unknown<->residual mappings to find non-zero partial derivatives and come up with a matrix-free equation for a single entry of the output, treating the derivatives as placeholders

<> We then use compile-time auto diff to generate an expression for each non-zero partial derivative, slotting them into the equation from step 1.

<> Finally, we parallelize across outputs of the matrix product. For large problems this easily saturates even high-end GPUs.

# EXTENDING THIS SIMPLE MODEL

---

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

In the paper, we show how we can expand this approach to handle more domains and solver techniques.

In particular, we

- <> Demonstrate how you can also write energies over mixed domains of meshes and images.
- <> Show how we can handle different solver variants of Gauss-Newton like Levenberg-Marquardt.
- <> Explore the tradeoff between using completely matrix-free products on one hand, and selectively pre-computing some parts of the matrix products on the other, to improve performance.



## EXTENDING THIS SIMPLE MODEL

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

1. Energies can be defined over **mixed domains** (meshes + images).

In the paper, we show how we can expand this approach to handle more domains and solver techniques.

In particular, we

- <> Demonstrate how you can also write energies over mixed domains of meshes and images.
- <> Show how we can handle different solver variants of Gauss-Newton like Levenberg-Marquardt.
- <> Explore the tradeoff between using completely matrix-free products on one hand, and selectively pre-computing some parts of the matrix products on the other, to improve performance.

## EXTENDING THIS SIMPLE MODEL

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

1. Energies can be defined over **mixed domains** (meshes + images).
2. Support of **Gauss-Newton variants**, like Levenberg-Marquardt.

In the paper, we show how we can expand this approach to handle more domains and solver techniques.

In particular, we

- <> Demonstrate how you can also write energies over mixed domains of meshes and images.
- <> Show how we can handle different solver variants of Gauss-Newton like Levenberg-Marquardt.
- <> Explore the tradeoff between using completely matrix-free products on one hand, and selectively pre-computing some parts of the matrix products on the other, to improve performance.

## EXTENDING THIS SIMPLE MODEL

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

1. Energies can be defined over **mixed domains** (meshes + images).
2. Support of **Gauss-Newton variants**, like Levenberg-Marquardt.
3. Primitives to tradeoff between completely matrix-free and **selectively precomputing** parts of the matrix expression before the inner PCG loop.

In the paper, we show how we can expand this approach to handle more domains and solver techniques.

In particular, we

<> Demonstrate how you can also write energies over mixed domains of meshes and images.

<> Show how we can handle different solver variants of Gauss-Newton like Levenberg-Marquardt.

<> Explore the tradeoff between using completely matrix-free products on one hand, and selectively pre-computing some parts of the matrix products on the other, to improve performance.



Now that we have a flavor for what Opt does, we can evaluate it along several axes

# OPT IS EXPRESSIVE

---

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

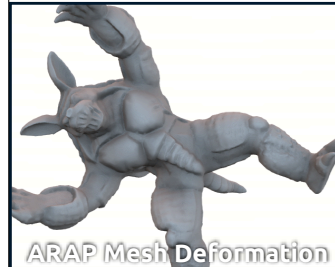
First, Opt is expressive enough to handle a wide range of problems. We implemented

- <> As-rigid-as-possible Mesh Deformation
- <> Poisson Image Editing
- <> Image Warping
- <> Shape from Shading
- <> Optical Flow
- <> Cotangent Mesh Smoothing
- <> Intrinsic Image Decomposition
- <> and Volumetric Mesh Deformation

all within Opt. This is only a subset of the applications we implemented for the paper. Each of these solvers only took

# OPT IS EXPRESSIVE

GENERATIONS / VANCOUVER  
SIGGRAPH 2018



© 2018 SIGGRAPH. All Rights Reserved

29

First, Opt is expressive enough to handle a wide range of problems. We implemented

- <> As-rigid-as-possible Mesh Deformation
- <> Poisson Image Editing
- <> Image Warping
- <> Shape from Shading
- <> Optical Flow
- <> Cotangent Mesh Smoothing
- <> Intrinsic Image Decomposition
- <> and Volumetric Mesh Deformation

all within Opt. This is only a subset of the applications we implemented for the paper. Each of these solvers only took

# OPT IS EXPRESSIVE

GENERATIONS / VANCOUVER  
SIGGRAPH 2018



First, Opt is expressive enough to handle a wide range of problems. We implemented

- <> As-rigid-as-possible Mesh Deformation
- <> Poisson Image Editing
- <> Image Warping
- <> Shape from Shading
- <> Optical Flow
- <> Cotangent Mesh Smoothing
- <> Intrinsic Image Decomposition
- <> and Volumetric Mesh Deformation

all within Opt. This is only a subset of the applications we implemented for the paper. Each of these solvers only took

# OPT IS EXPRESSIVE

GENERATIONS / VANCOUVER  
SIGGRAPH 2018



© 2018 SIGGRAPH. All Rights Reserved

29

First, Opt is expressive enough to handle a wide range of problems. We implemented

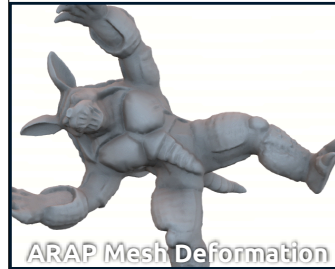
- <> As-rigid-as-possible Mesh Deformation
- <> Poisson Image Editing
- <> Image Warping
- <> Shape from Shading
- <> Optical Flow
- <> Cotangent Mesh Smoothing
- <> Intrinsic Image Decomposition
- <> and Volumetric Mesh Deformation

all within Opt. This is only a subset of the applications we implemented for the paper. Each of these solvers only took



# OPT IS EXPRESSIVE

GENERATIONS / VANCOUVER  
SIGGRAPH 2018



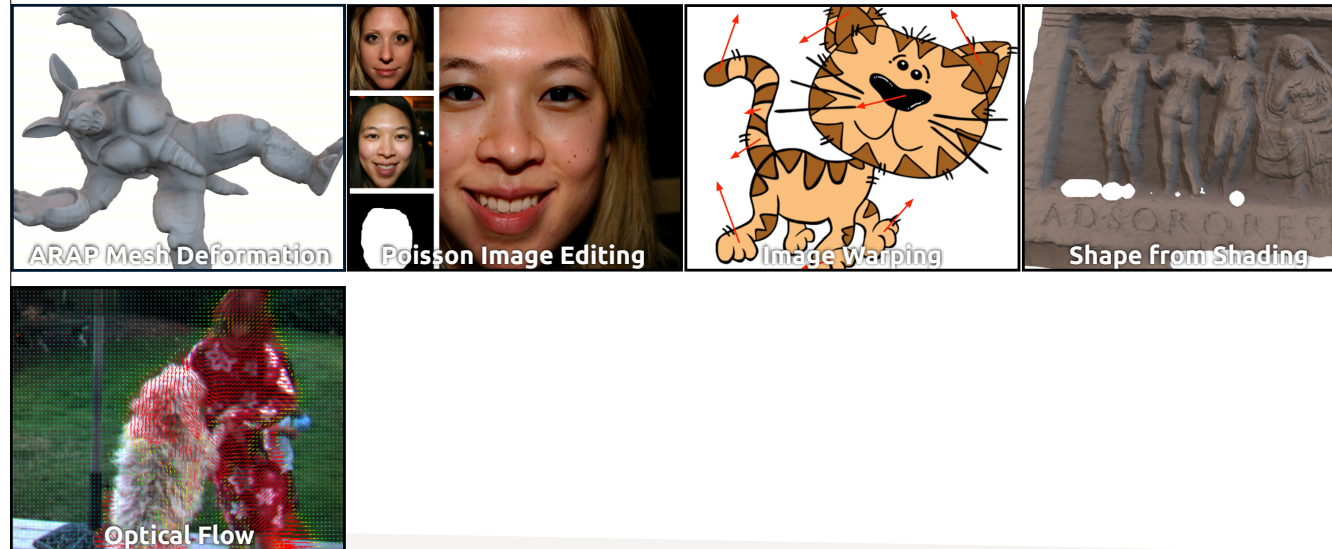
First, Opt is expressive enough to handle a wide range of problems. We implemented

- <> As-rigid-as-possible Mesh Deformation
- <> Poisson Image Editing
- <> Image Warping
- <> Shape from Shading
- <> Optical Flow
- <> Cotangent Mesh Smoothing
- <> Intrinsic Image Decomposition
- <> and Volumetric Mesh Deformation

all within Opt. This is only a subset of the applications we implemented for the paper. Each of these solvers only took

# OPT IS EXPRESSIVE

GENERATIONS / VANCOUVER  
SIGGRAPH 2018



© 2018 SIGGRAPH. All Rights Reserved

29

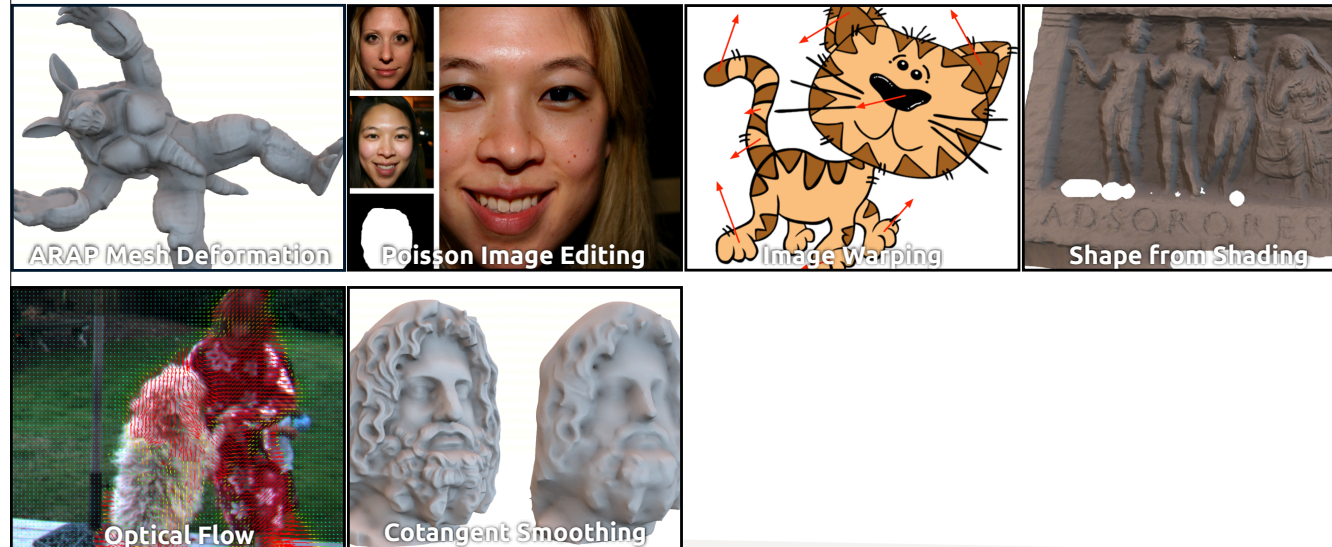
First, Opt is expressive enough to handle a wide range of problems. We implemented

- <> As-rigid-as-possible Mesh Deformation
- <> Poisson Image Editing
- <> Image Warping
- <> Shape from Shading
- <> Optical Flow
- <> Cotangent Mesh Smoothing
- <> Intrinsic Image Decomposition
- <> and Volumetric Mesh Deformation

all within Opt. This is only a subset of the applications we implemented for the paper. Each of these solvers only took

# OPT IS EXPRESSIVE

GENERATIONS / VANCOUVER  
SIGGRAPH 2018



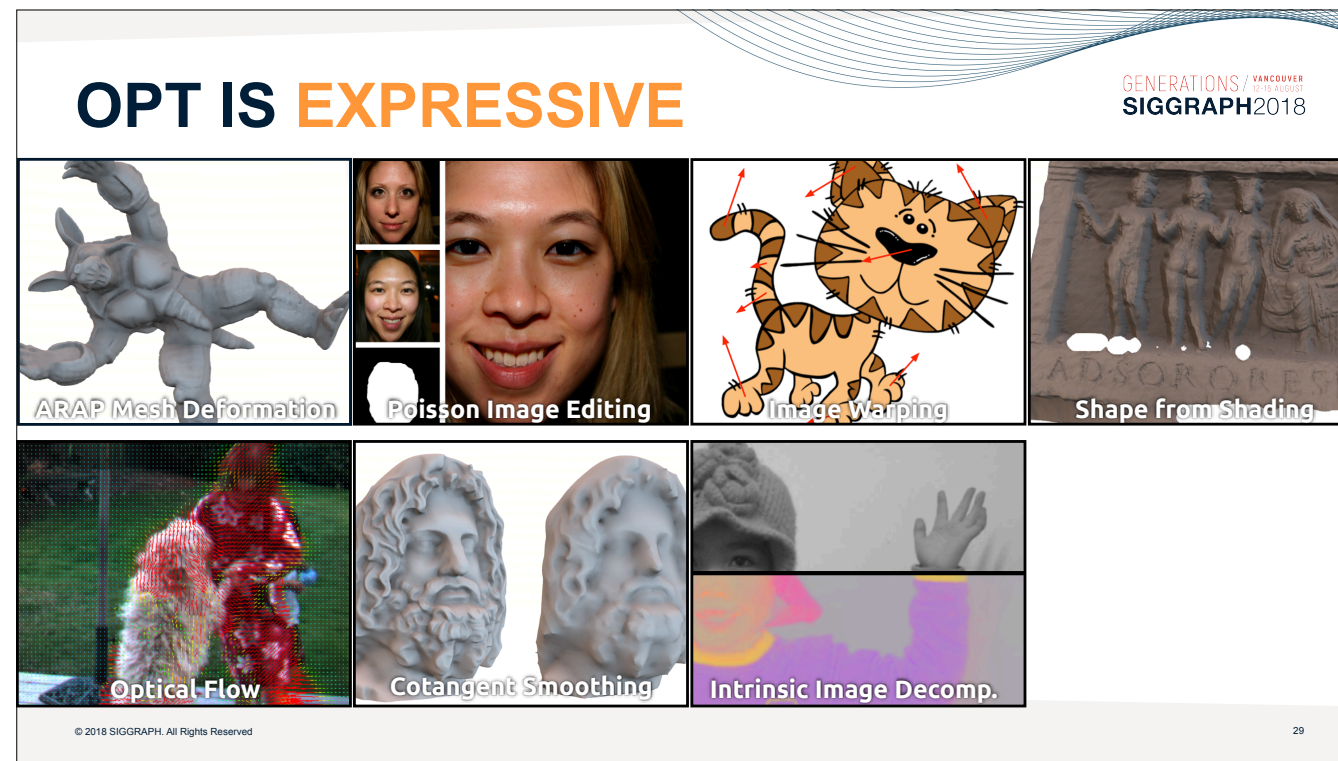
© 2018 SIGGRAPH. All Rights Reserved

29

First, Opt is expressive enough to handle a wide range of problems. We implemented

- <> As-rigid-as-possible Mesh Deformation
- <> Poisson Image Editing
- <> Image Warping
- <> Shape from Shading
- <> Optical Flow
- <> Cotangent Mesh Smoothing
- <> Intrinsic Image Decomposition
- <> and Volumetric Mesh Deformation

all within Opt. This is only a subset of the applications we implemented for the paper. Each of these solvers only took

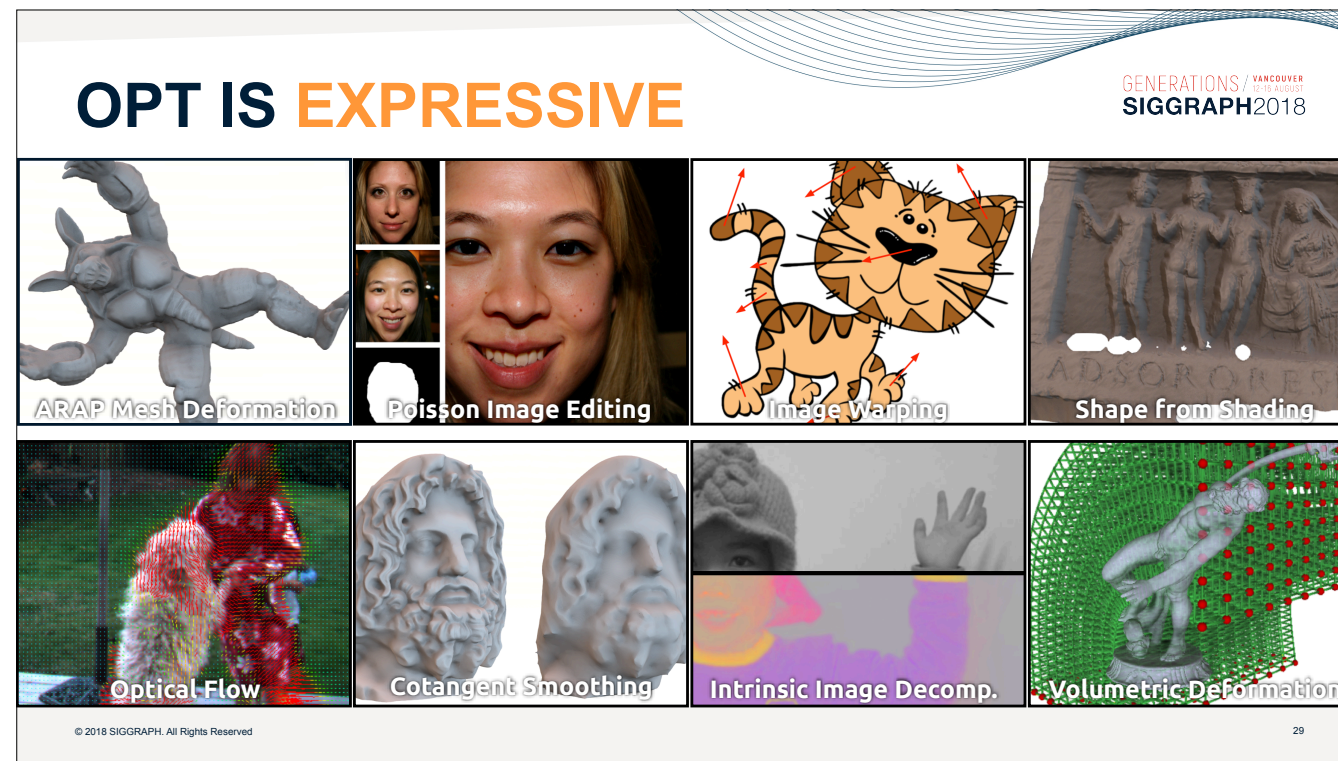


First, Opt is expressive enough to handle a wide range of problems. We implemented

- <> As-rigid-as-possible Mesh Deformation
- <> Poisson Image Editing
- <> Image Warping
- <> Shape from Shading
- <> Optical Flow
- <> Cotangent Mesh Smoothing
- <> Intrinsic Image Decomposition
- <> and Volumetric Mesh Deformation

all within Opt. This is only a subset of the applications we implemented for the paper. Each of these solvers only took









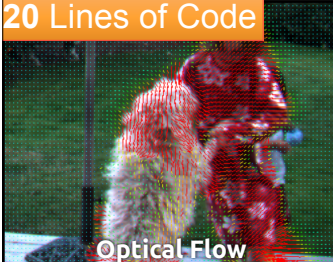
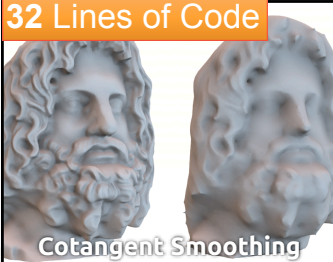

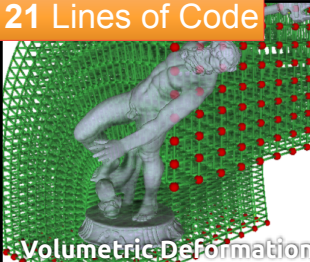
First, Opt is expressive enough to handle a wide range of problems. We implemented

- <> As-rigid-as-possible Mesh Deformation
- <> Poisson Image Editing
- <> Image Warping
- <> Shape from Shading
- <> Optical Flow
- <> Cotangent Mesh Smoothing
- <> Intrinsic Image Decomposition
- <> and Volumetric Mesh Deformation

all within Opt. This is only a subset of the applications we implemented for the paper. Each of these solvers only took

# OPT IS CONCISE

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

<b>18 Lines of Code</b>  ARAP Mesh Deformation	<b>13 Lines of Code</b>  Poisson Image Editing	<b>21 Lines of Code</b>  Image Warping	<b>96 Lines of Code</b>  Shape from Shading
<b>20 Lines of Code</b>  Optical Flow	<b>32 Lines of Code</b>  Cotangent Smoothing	<b>32 Lines of Code</b>  Intrinsic Image Decomp.	<b>21 Lines of Code</b>  Volumetric Deformation

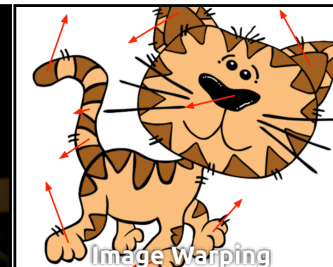
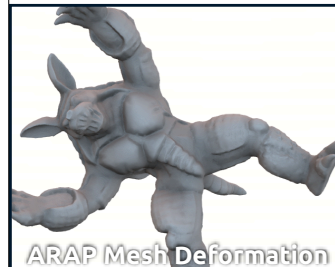
© 2018 SIGGRAPH. All Rights Reserved

30

a handful of lines of code to implement in Opt, indeed all but one were less than 40 lines.

# OPT IS CONCISE

GENERATIONS / VANCOUVER  
SIGGRAPH 2018



© 2018 SIGGRAPH. All Rights Reserved

Relative Lines of Code

31

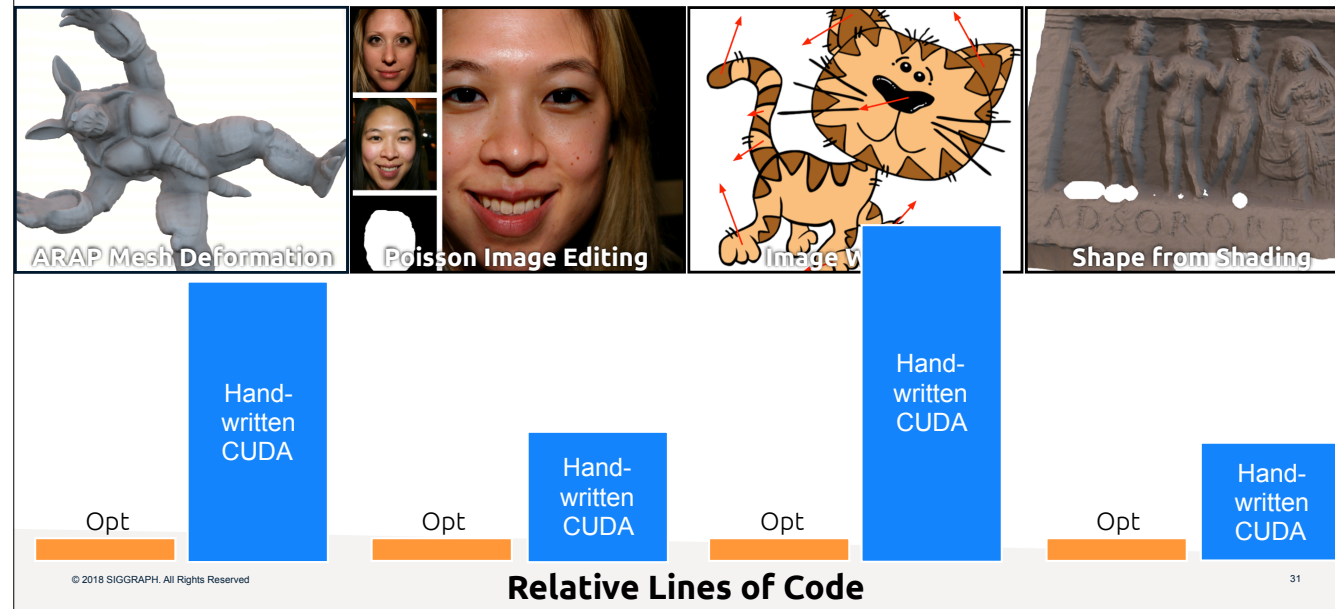
Four of these were previously implemented laboriously by hand in CUDA;

<> for these we can directly compare solver length and see Opt code is far more compact. Every solver is at least 4.5x more verbose in CUDA, and the worst is over 13x longer!

Note that an additional energy term requires adding a couple of lines of Opt code, but in the handwritten solvers requires surgery on at least three different pieces of code, massively increasing chances for mismatch errors or problems with by-hand differentiation.

# OPT IS CONCISE

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

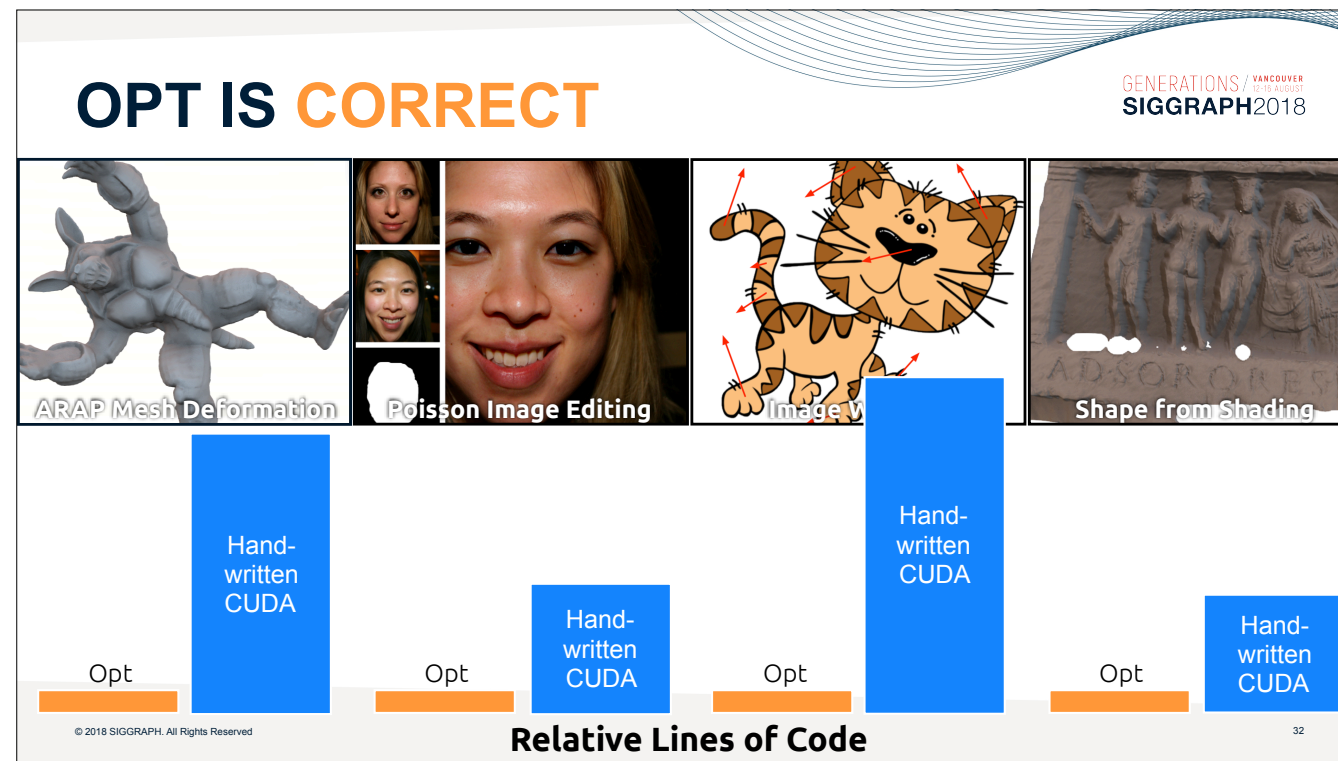


Four of these were previously implemented laboriously by hand in CUDA;

<> for these we can directly compare solver length and see Opt code is far more compact. Every solver is at least 4.5x more verbose in CUDA, and the worst is over 13x longer!

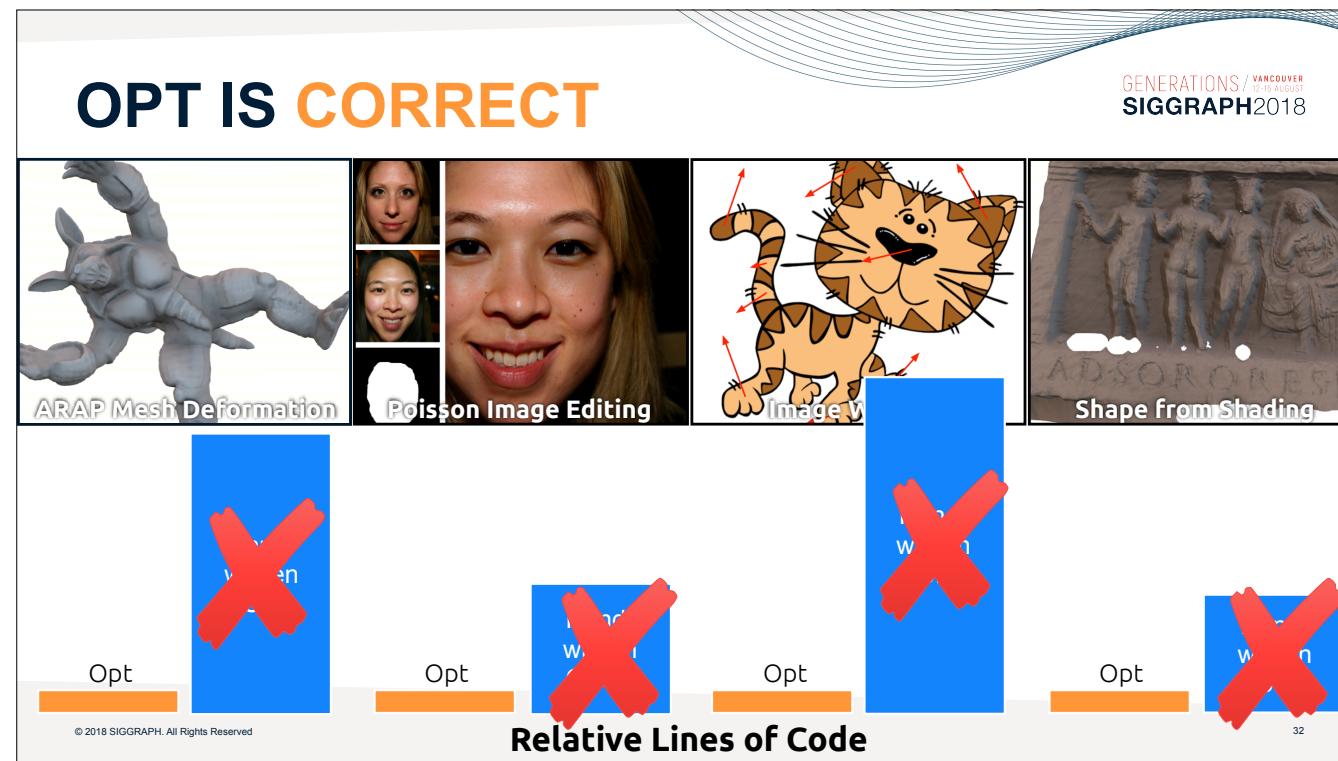
Note that an additional energy term requires adding a couple of lines of Opt code, but in the handwritten solvers requires surgery on at least three different pieces of code, massively increasing chances for mismatch errors or problems with by-hand differentiation.





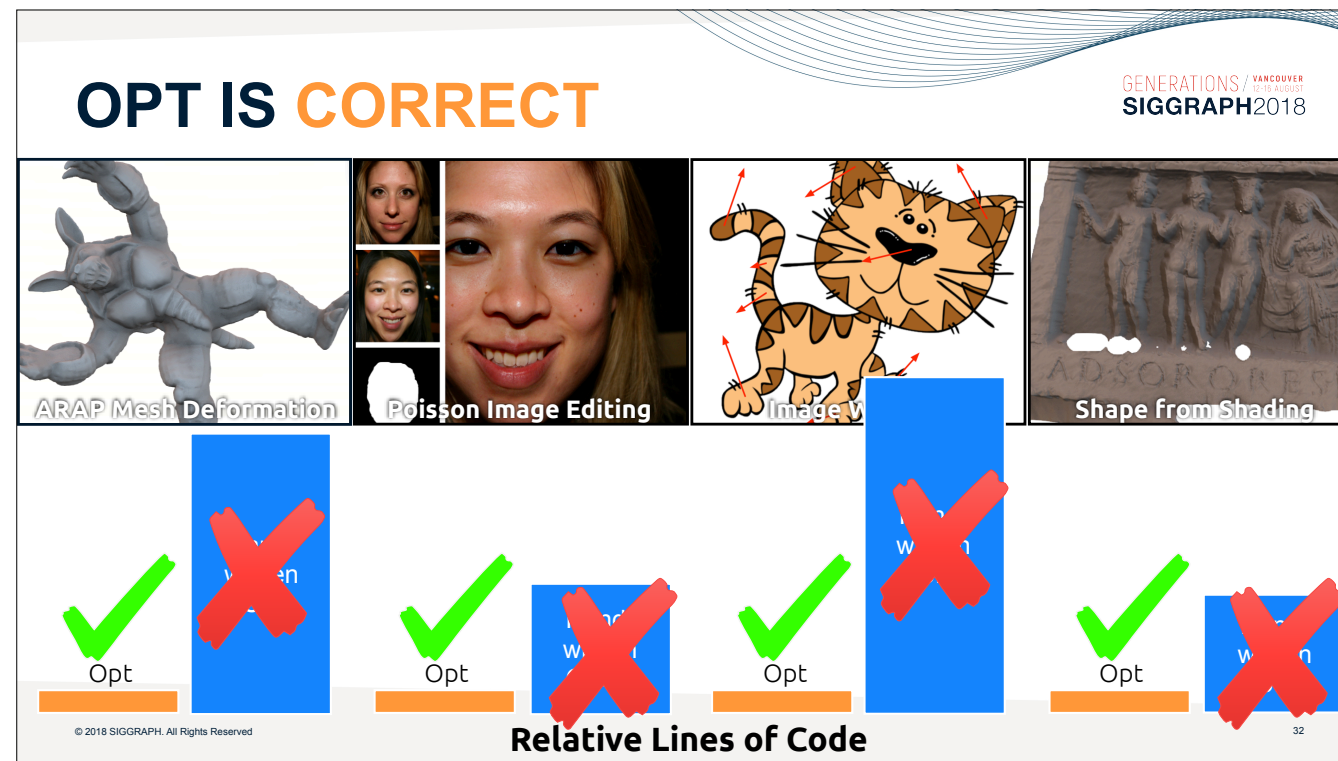
In fact, every single handwritten solver we compared against had at least some error in the derivative terms, either in the calculus or boundary conditions, which negatively impacted solver convergence until we fixed them.

<> Opt, by offloading the differentiation labor and code correspondence bookkeeping to the compiler, generates solvers that are correct by construction.



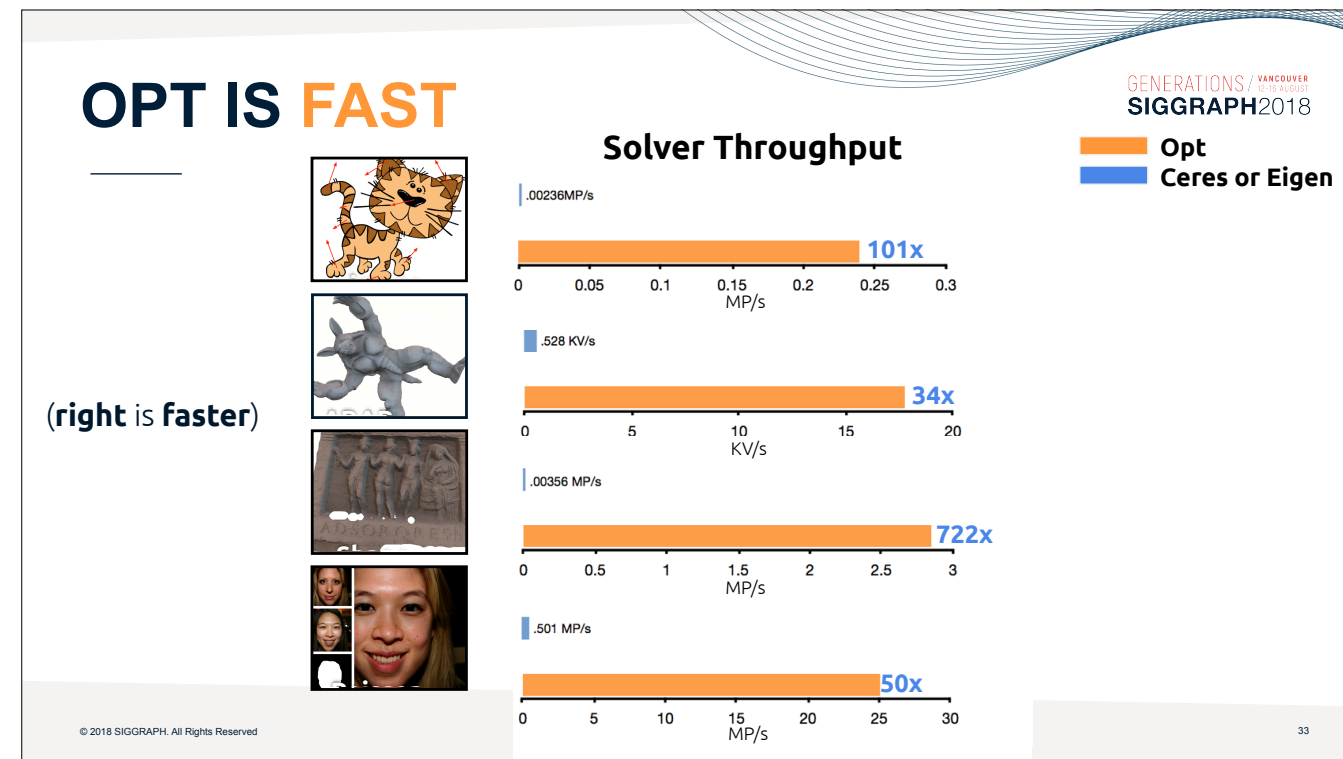
In fact, every single handwritten solver we compared against had at least some error in the derivative terms, either in the calculus or boundary conditions, which negatively impacted solver convergence until we fixed them.

<> Opt, by offloading the differentiation labor and code correspondence bookkeeping to the compiler, generates solvers that are correct by construction.



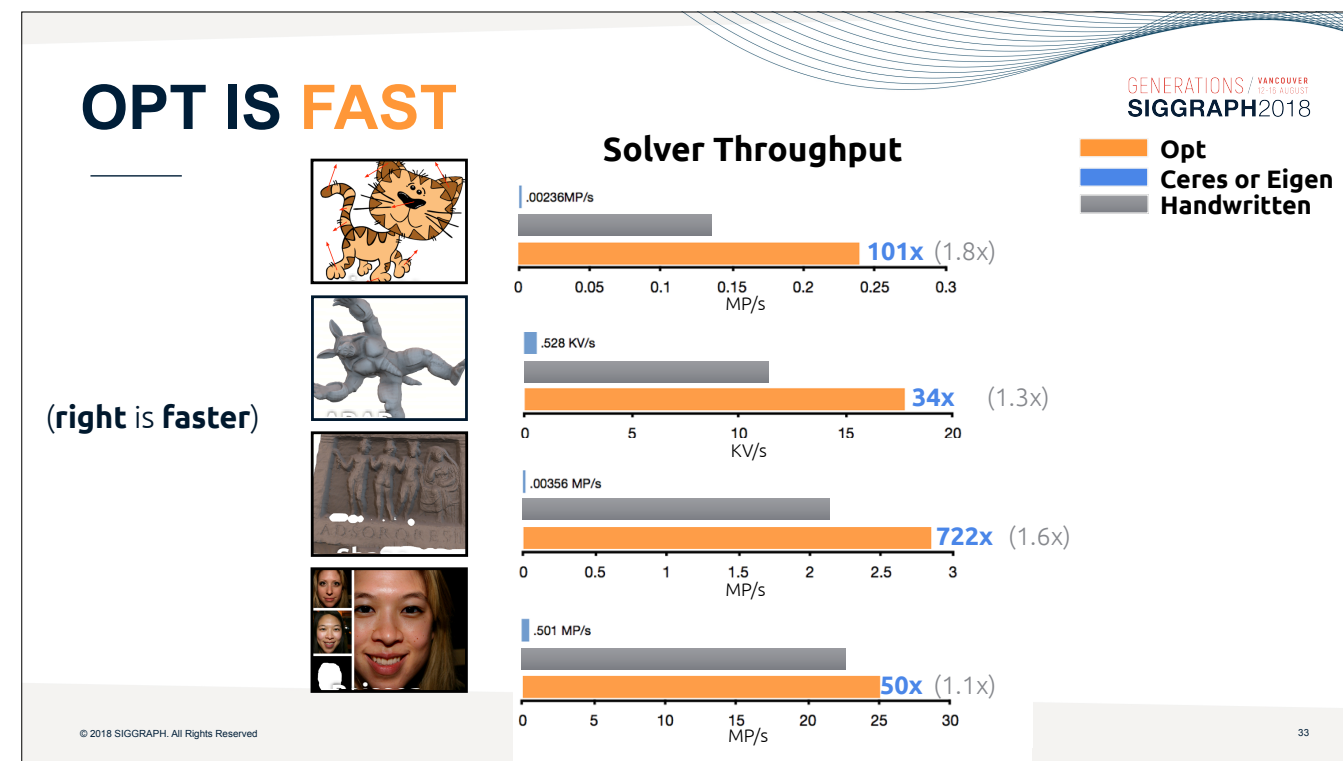
In fact, every single handwritten solver we compared against had at least some error in the derivative terms, either in the calculus or boundary conditions, which negatively impacted solver convergence until we fixed them.

<> Opt, by offloading the differentiation labor and code correspondence bookkeeping to the compiler, generates solvers that are correct by construction.



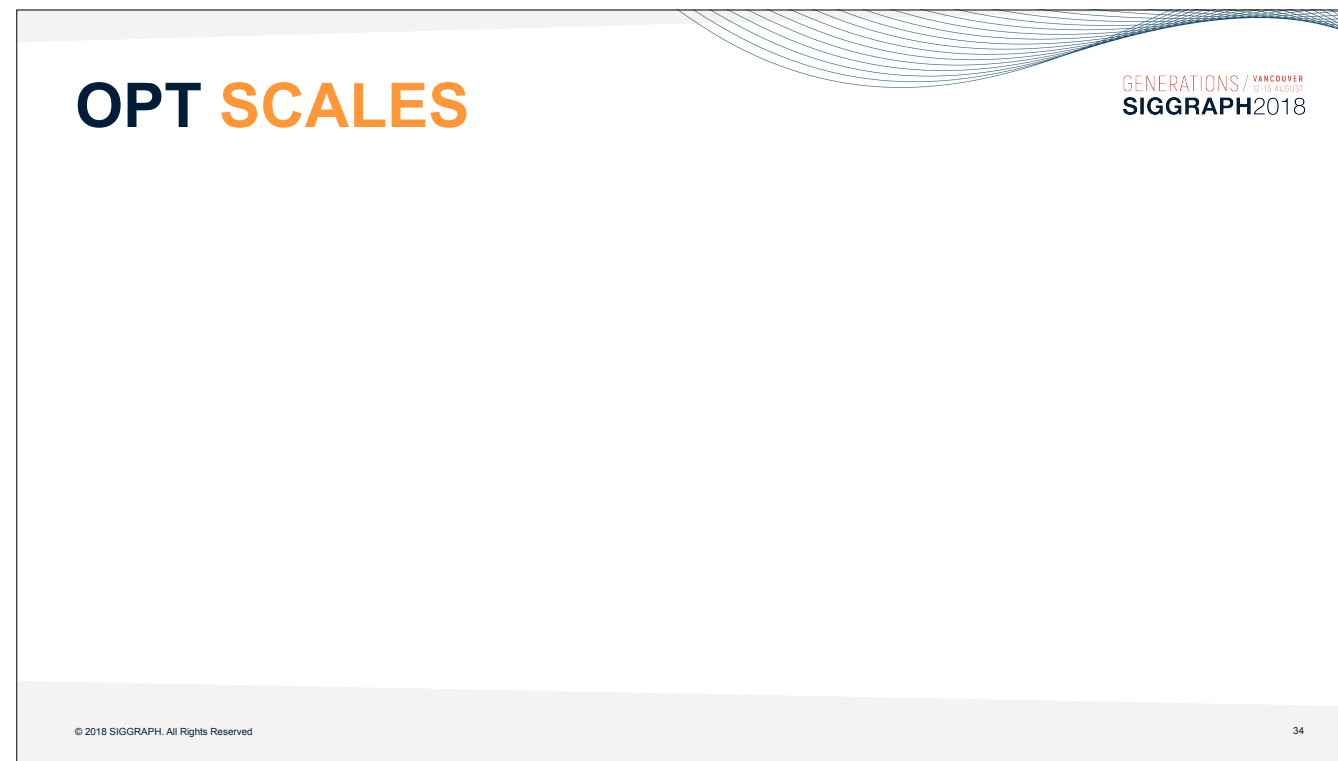
And performance wise, Opt blows away other high level solvers, by multiple orders of magnitude. Here we have the throughput of solvers implemented in Opt (shown in orange) and Ceres (or Eigen in the case of Poisson Image Editing), shown in blue. The lowest speedup we get with Opt is 33x on Mesh Deformation, the highest is over 700 times on our shape from shading implementation.

<> It even beats out all of the handwritten solvers that we compared against (shown in gray), ranging from 1.1x to 1.8x faster. Our compiler does some algebraic simplifications and boundary handling that the original authors of the handwritten implementations either did notice or did not bother with, you can read more about these in the paper.



And performance wise, Opt blows away other high level solvers, by multiple orders of magnitude. Here we have the throughput of solvers implemented in Opt (shown in orange) and Ceres (or Eigen in the case of Poisson Image Editing), shown in blue. The lowest speedup we get with Opt is 33x on Mesh Deformation, the highest is over 700 times on our shape from shading implementation.

<> It even beats out all of the handwritten solvers that we compared against (shown in gray), ranging from 1.1x to 1.8x faster. Our compiler does some algebraic simplifications and boundary handling that the original authors of the handwritten implementations either did notice or did not bother with, you can read more about these in the paper.

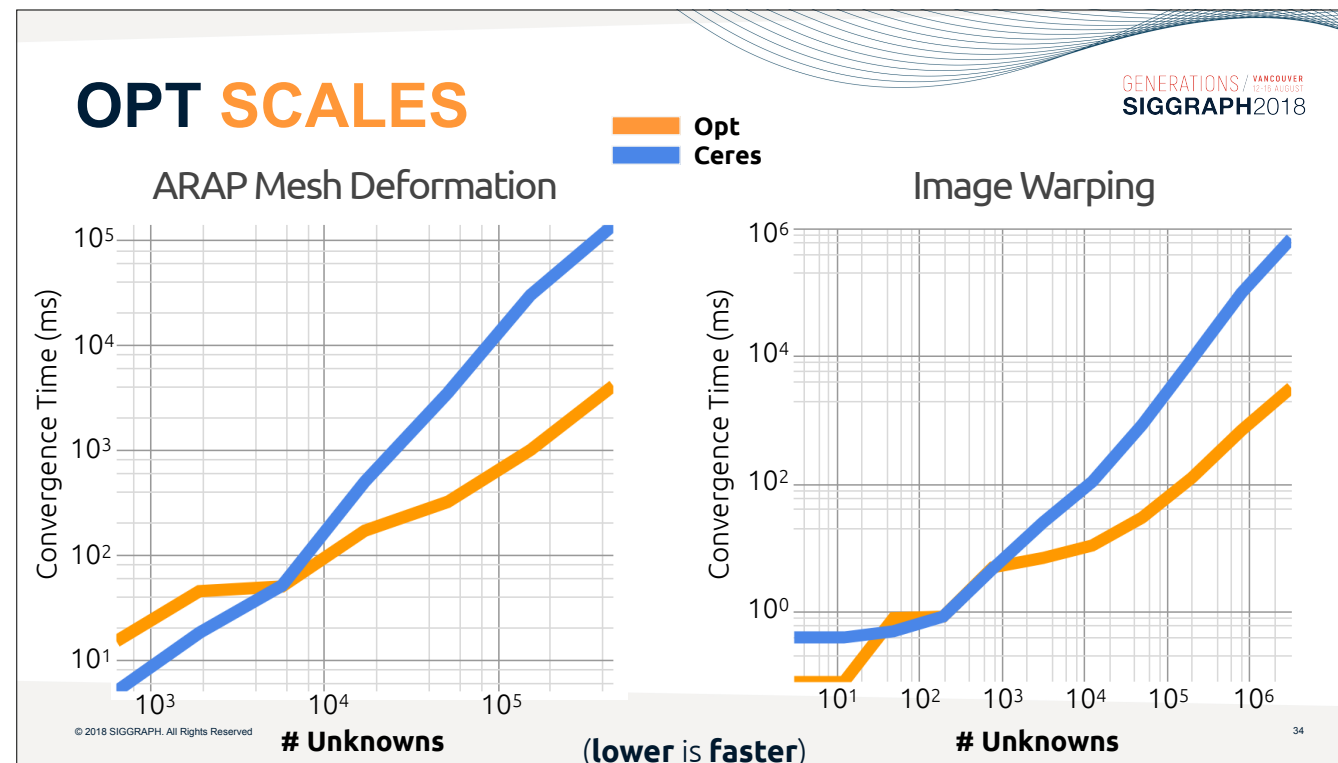


And here we see the performance of Opt versus Ceres as we vary problem size. <>

Again, Opt is in orange and Ceres is in Blue. We are charting Convergence time vs # of unknowns, on a log-log chart, so large differences are quite compressed.

At low unknown count (in the several hundreds) high-level CPU solvers that do not have to transfer data back and forth from the GPU co-processor compare favorably to Opt,

<> but once we get into larger problems with thousands unknowns or more, where the massive parallelism of the GPU comes into play, Opt wins out in performance, and as size increases, so does the performance gap.

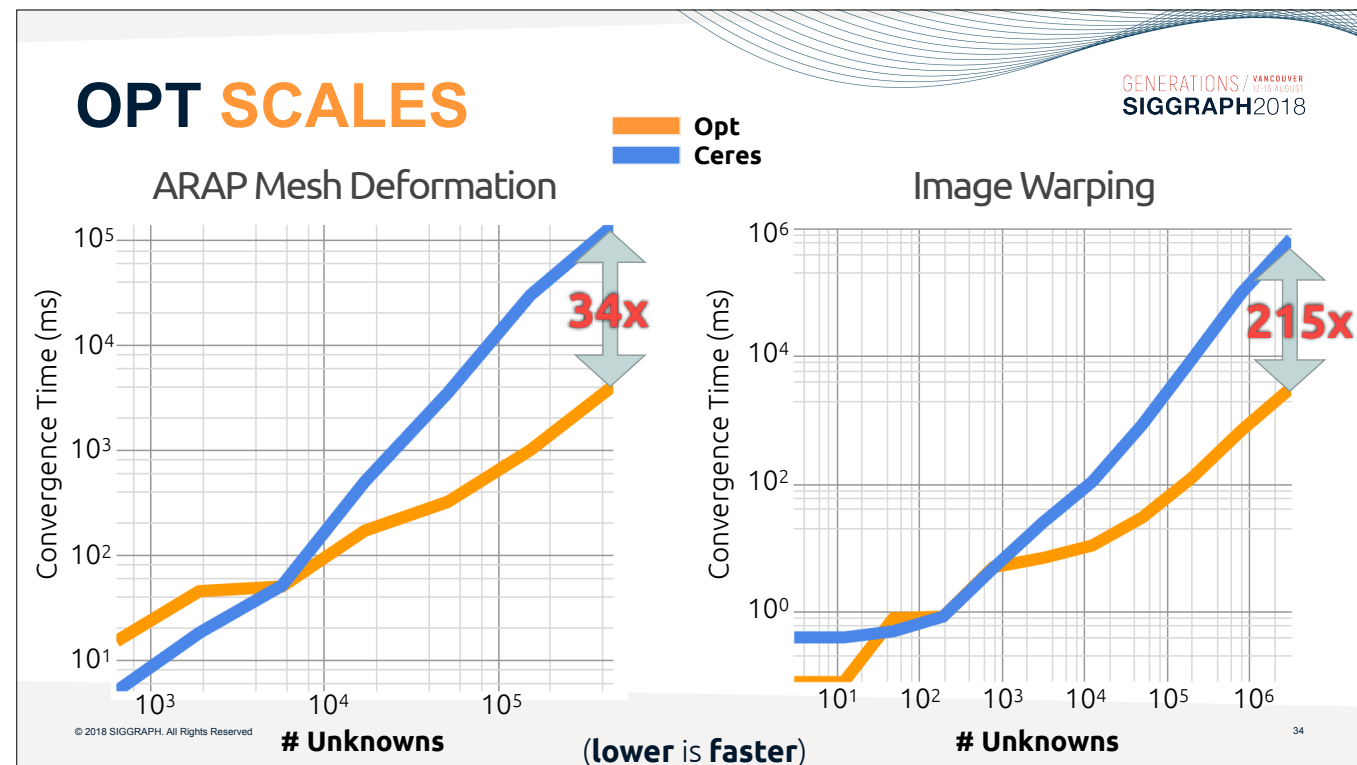


And here we see the performance of Opt versus Ceres as we vary problem size. <>

Again, Opt is in orange and Ceres is in Blue. We are charting Convergence time vs # of unknowns, on a log-log chart, so large differences are quite compressed.

At low unknown count (in the several hundreds) high-level CPU solvers that do not have to transfer data back and forth from the GPU co-processor compare favorably to Opt,

<> but once we get into larger problems with thousands unknowns or more, where the massive parallelism of the GPU comes into play, Opt wins out in performance, and as size increases, so does the performance gap.



And here we see the performance of Opt versus Ceres as we vary problem size. <>

Again, Opt is in orange and Ceres is in Blue. We are charting Convergence time vs # of unknowns, on a log-log chart, so large differences are quite compressed.

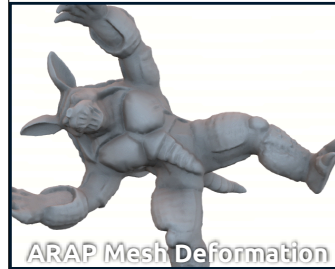
At low unknown count (in the several hundreds) high-level CPU solvers that do not have to transfer data back and forth from the GPU co-processor compare favorably to Opt,

<> but once we get into larger problems with thousands unknowns or more, where the massive parallelism of the GPU comes into play, Opt wins out in performance, and as size increases, so does the performance gap.



## MATRIX-FREE IS IMPORTANT FOR SPEED

GENERATIONS / VANCOUVER  
SIGGRAPH 2018



ARAP Mesh Deformation



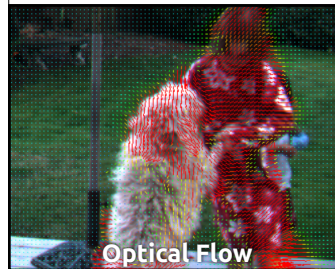
Poisson Image Editing



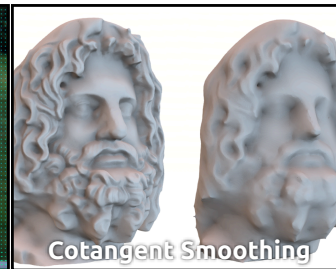
Image Warping



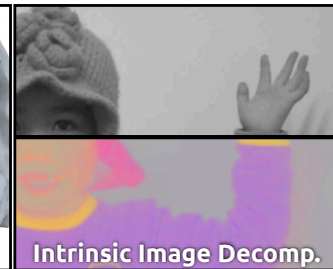
Shape from Shading



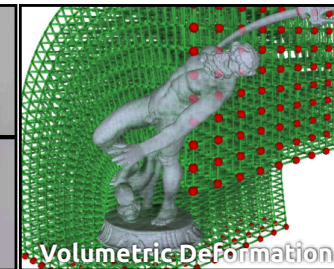
Optical Flow



Cotangent Smoothing



Intrinsic Image Decomposition



Volumetric Deformation

## MATRIX-FREE IS IMPORTANT FOR SPEED

GENERATIONS / VANCOUVER  
SIGGRAPH 2018



© 2018 SIGGRAPH. All Rights Reserved

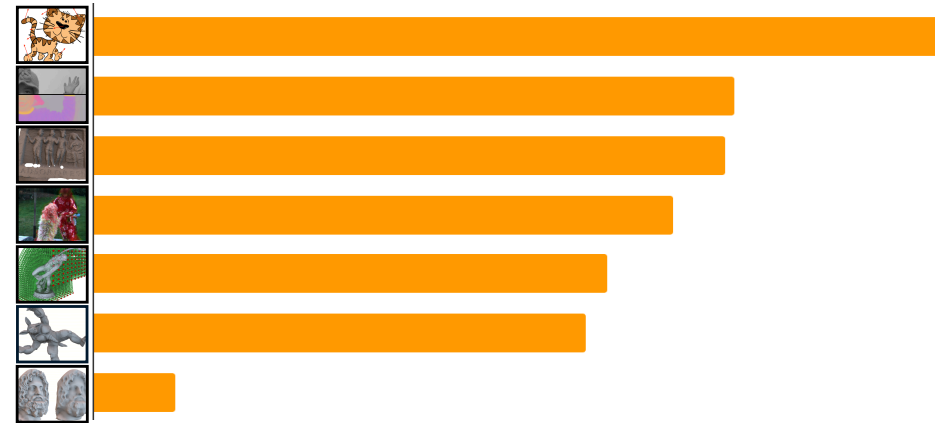
36

The ability for Opt to generate Matrix-Free code is important to get improved performance. Here we chart the

<> relative performance of our examples for doing PCG iterations using Opt-generated matrix free code, using PCG with a materialized JtJ matrix (in compressed-row storage form) as the baseline. For all but one of our problems we get a significant performance improvement by choosing the matrix free approach. And this chart doesn't reflect the fact that materializing the matrix in the first place has a fixed overhead that our matrix-free code does not pay.

## MATRIX-FREE IS IMPORTANT FOR SPEED

GENERATIONS / VANCOUVER  
SIGGRAPH 2018



Relative Performance of Matrix-Free PCG Iterations (to the right is faster)

© 2018 SIGGRAPH. All Rights Reserved

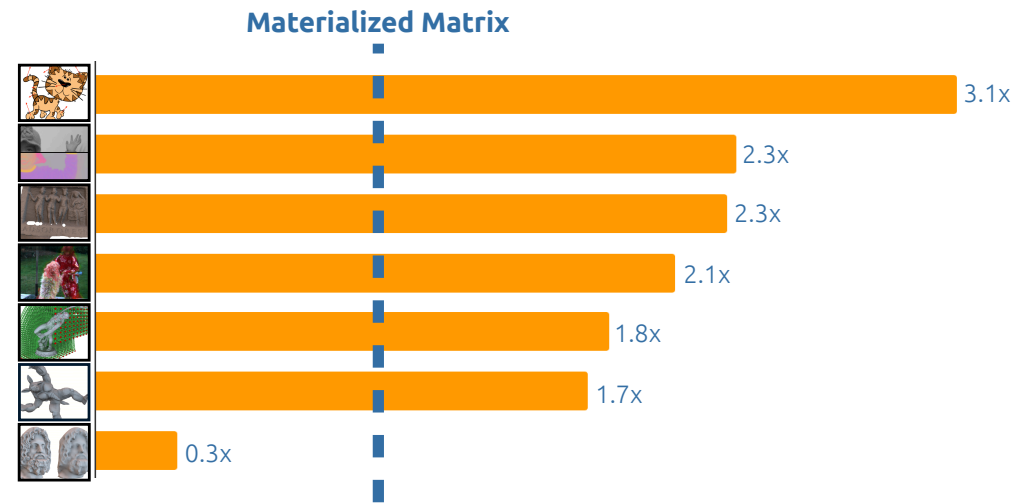
36

The ability for Opt to generate Matrix-Free code is important to get improved performance. Here we chart the

<> relative performance of our examples for doing PCG iterations using Opt-generated matrix free code, using PCG with a materialized JtJ matrix (in compressed-row storage form) as the baseline. For all but one of our problems we get a significant performance improvement by choosing the matrix free approach. And this chart doesn't reflect the fact that materializing the matrix in the first place has a fixed overhead that our matrix-free code does not pay.

## MATRIX-FREE IS IMPORTANT FOR SPEED

GENERATIONS / VANCOUVER  
SIGGRAPH 2018



Relative Performance of Matrix-Free PCG Iterations (to the right is faster)

© 2018 SIGGRAPH. All Rights Reserved

36

The ability for Opt to generate Matrix-Free code is important to get improved performance. Here we chart the

<> relative performance of our examples for doing PCG iterations using Opt-generated matrix free code, using PCG with a materialized JtJ matrix (in compressed-row storage form) as the baseline. For all but one of our problems we get a significant performance improvement by choosing the matrix free approach. And this chart doesn't reflect the fact that materializing the matrix in the first place has a fixed overhead that our matrix-free code does not pay.

## WHAT MAKES OPT FAST?

---

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

So ultimately, what makes Opt fast?

<> Well first, because of our restricted DSL and exploiting structure in the energy, we can generate code for the Jacobian that can take advantage of massive parallelism and wide-SIMD units, like on the GPU. This is responsible for a large factor of the speed difference between Opt and other high-level solvers

<> Second, we can go a step further and remove the potential overhead of materializing an explicit system matrix, which we just saw often leads to further performance improvements

## WHAT MAKES OPT FAST?

GENERATIONS / VANCOUVER  
SIGGRAPH2018

- GPU-ifies J construction

So ultimately, what makes Opt fast?

<> Well first, because of our restricted DSL and exploiting structure in the energy, we can generate code for the Jacobian that can take advantage of massive parallelism and wide-SIMD units, like on the GPU. This is responsible for a large factor of the speed difference between Opt and other high-level solvers

<> Second, we can go a step further and remove the potential overhead of materializing an explicit system matrix, which we just saw often leads to further performance improvements

## WHAT MAKES OPT FAST?

GENERATIONS / VANCOUVER  
SIGGRAPH2018

- GPU-ifies J construction
- Removes explicit matrix

So ultimately, what makes Opt fast?

<> Well first, because of our restricted DSL and exploiting structure in the energy, we can generate code for the Jacobian that can take advantage of massive parallelism and wide-SIMD units, like on the GPU. This is responsible for a large factor of the speed difference between Opt and other high-level solvers

<> Second, we can go a step further and remove the potential overhead of materializing an explicit system matrix, which we just saw often leads to further performance improvements



Our design decisions provide us with some natural limitations

<> which suggests future work.

<> Our access patterns are restricted, which is the very thing ensuring we can do the data dependency analysis we need to to maximize parallelism. Future work would relax these restrictions

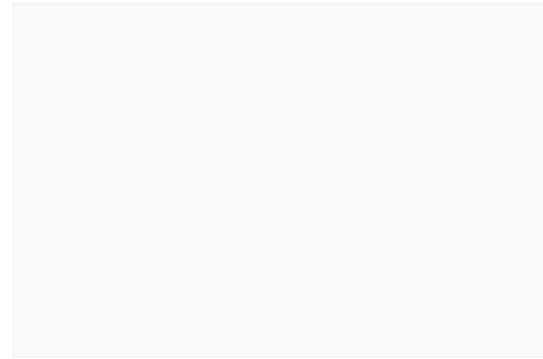
<> We allow for selectively precomputing parts of the Jacobian computation, which allows users to eke out extra performance, but it should be possible to generalize to a much wider range of schedules, like Halide.

<> We have fast matrix-free code generation for a small set of matrix calculus primitives, those necessary to generate Gauss-Newton-like solvers. Having a language that allows arbitrary tensor derivatives would open up a large set of new potential use-cases that are currently too laborious to explore.

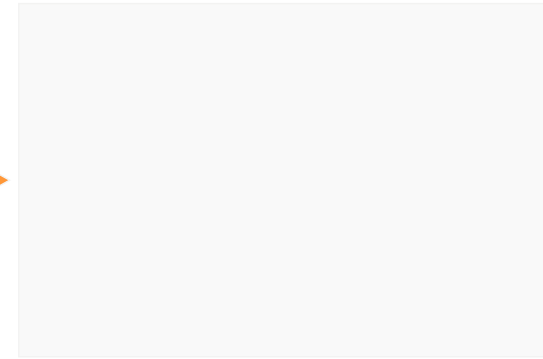


# LIMITATIONS AND FUTURE WORK

GENERATIONS / VANCOUVER  
SIGGRAPH 2018



**Limitations**



**Future Work**

© 2018 SIGGRAPH. All Rights Reserved

38

Our design decisions provide us with some natural limitations

<> which suggests future work.

<> Our access patterns are restricted, which is the very thing ensuring we can do the data dependency analysis we need to to maximize parallelism. Future work would relax these restrictions

<> We allow for selectively precomputing parts of the Jacobian computation, which allows users to eke out extra performance, but it should be possible to generalize to a much wider range of schedules, like Halide.

<> We have fast matrix-free code generation for a small set of matrix calculus primitives, those necessary to generate Gauss-Newton-like solvers. Having a language that allows arbitrary tensor derivatives would open up a large set of new potential use-cases that are currently too laborious to explore.

# LIMITATIONS AND FUTURE WORK

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

- Restricted access patterns



- More access patterns

**Limitations**

**Future Work**

© 2018 SIGGRAPH. All Rights Reserved

38

Our design decisions provide us with some natural limitations

<> which suggests future work.

<> Our access patterns are restricted, which is the very thing ensuring we can do the data dependency analysis we need to to maximize parallelism. Future work would relax these restrictions

<> We allow for selectively precomputing parts of the Jacobian computation, which allows users to eke out extra performance, but it should be possible to generalize to a much wider range of schedules, like Halide.

<> We have fast matrix-free code generation for a small set of matrix calculus primitives, those necessary to generate Gauss-Newton-like solvers. Having a language that allows arbitrary tensor derivatives would open up a large set of new potential use-cases that are currently too laborious to explore.

## LIMITATIONS AND FUTURE WORK

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

- Restricted access patterns
- Limited scheduling



- More access patterns
- Halide-like scheduling

**Limitations**

**Future Work**

© 2018 SIGGRAPH. All Rights Reserved

38

Our design decisions provide us with some natural limitations

<> which suggests future work.

<> Our access patterns are restricted, which is the very thing ensuring we can do the data dependency analysis we need to to maximize parallelism. Future work would relax these restrictions

<> We allow for selectively precomputing parts of the Jacobian computation, which allows users to eke out extra performance, but it should be possible to generalize to a much wider range of schedules, like Halide.

<> We have fast matrix-free code generation for a small set of matrix calculus primitives, those necessary to generate Gauss-Newton-like solvers. Having a language that allows arbitrary tensor derivatives would open up a large set of new potential use-cases that are currently too laborious to explore.

## LIMITATIONS AND FUTURE WORK

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

- Restricted access patterns
- Limited scheduling
- Only GN-like solvers



- More access patterns
- Halide-like scheduling
- General matrix calculus terms

**Limitations**

**Future Work**

© 2018 SIGGRAPH. All Rights Reserved

38

Our design decisions provide us with some natural limitations

<> which suggests future work.

<> Our access patterns are restricted, which is the very thing ensuring we can do the data dependency analysis we need to to maximize parallelism. Future work would relax these restrictions

<> We allow for selectively precomputing parts of the Jacobian computation, which allows users to eke out extra performance, but it should be possible to generalize to a much wider range of schedules, like Halide.

<> We have fast matrix-free code generation for a small set of matrix calculus primitives, those necessary to generate Gauss-Newton-like solvers. Having a language that allows arbitrary tensor derivatives would open up a large set of new potential use-cases that are currently too laborious to explore.

# WRAP-UP

---

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

To wrap things up,

- <> we solved the problem we set out to: Very fast non-linear least squares optimizers on images/meshes/graphs are easy to write
- <> We have several actually working implementations of real problems including recent SIGGRAPH papers
- <> And hundreds of people have used it in one form or another
- <> Opt has its own website, and an active community on github, where its released under the open-source MIT license.

Thank you for your attention.

# WRAP-UP

---

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

- Very fast NLLS optimizers on images/meshes/graphs are easy to write

To wrap things up,

- <> we solved the problem we set out to: Very fast non-linear least squares optimizers on images/meshes/graphs are easy to write
- <> We have several actually working implementations of real problems including recent SIGGRAPH papers
- <> And hundreds of people have used it in one form or another

<> Opt has its own website, and an active community on github, where its released under the open-source MIT license.

Thank you for your attention.

# WRAP-UP

---

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

- Very fast NLLS optimizers on images/meshes/graphs are easy to write
- 10+ working examples

To wrap things up,

- <> we solved the problem we set out to: Very fast non-linear least squares optimizers on images/meshes/graphs are easy to write
- <> We have several actually working implementations of real problems including recent SIGGRAPH papers
- <> And hundreds of people have used it in one form or another

<> Opt has its own website, and an active community on github, where its released under the open-source MIT license.

Thank you for your attention.

# WRAP-UP

---

GENERATIONS / VANCOUVER  
SIGGRAPH 2018

- Very fast NLLS optimizers on images/meshes/graphs are easy to write
- 10+ working examples
- Many people use it

To wrap things up,

- <> we solved the problem we set out to: Very fast non-linear least squares optimizers on images/meshes/graphs are easy to write
- <> We have several actually working implementations of real problems including recent SIGGRAPH papers
- <> And hundreds of people have used it in one form or another

<> Opt has its own website, and an active community on github, where its released under the open-source MIT license.

Thank you for your attention.



# WRAP-UP

GENERATIONS / VANCOUVER  
SIGGRAPH2018

- Very fast NLLS optimizers on images/meshes/graphs are easy to write
- 10+ working examples
- Many people use it

[optlang.org/](http://optlang.org/)

Open source under the MIT license



To wrap things up,

- <> we solved the problem we set out to: Very fast non-linear least squares optimizers on images/meshes/graphs are easy to write
- <> We have several actually working implementations of real problems including recent SIGGRAPH papers
- <> And hundreds of people have used it in one form or another

<> Opt has its own website, and an active community on github, where its released under the open-source MIT license.

Thank you for your attention.

# THANK YOU!

---

GENERATIONS / VANCOUVER  
SIGGRAPH2018

[optlang.org/](http://optlang.org/)  
Open source under the MIT license





GENERATIONS / VANCOUVER  
12-16 AUGUST  
**SIGGRAPH2018**

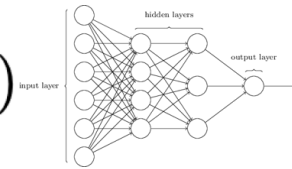
## EXTRA SLIDES



# Why Gauss-Newton?

Gradient Descent?

$$\mathbf{a}_{n+1} = \mathbf{a}_n - \gamma \nabla F(\mathbf{a}_n)$$



Newton's Method

$$x_{n+1} = x_n - \frac{f'(x_n)}{f''(x_n)}$$

Single variable...

$$\mathbf{x}_{n+1} = \mathbf{x}_n - [\mathbf{H}f(\mathbf{x}_n)]^{-1} \nabla f(\mathbf{x}_n)$$

$$[\mathbf{H}f(\mathbf{x}_n)] \Delta \mathbf{x} = -\nabla f(\mathbf{x}_n)$$

Why are we doing Gauss-Newton? (and what is it anyway?)

Why don't we use gradient descent for optimization like all the cool kids? It's simple to understand, and obviously works.

Well, it actually has bad convergence on interesting functions. There is a reason even beyond sheer data size why neural nets take forever to train.

But simplicity is good! The problem with gradient descent is we aren't using much information, we are linearly approximating our function at each step to find our next move. If we quadratically approximate instead, we get Newton's method for optimization, which is taught as a root finding method in high schools all over. That seems promisingly simple, while having better convergence properties. Of course what I put up there has a problem, it's one dimensional.

We can move to higher dimensions through analogy. The derivative becomes the gradient vector, and the second derivative is the Hessian matrix.

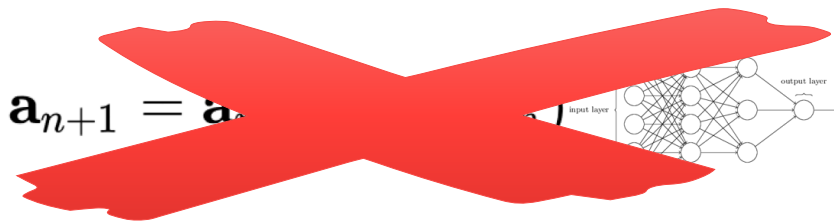
Problem, we can't divide by a matrix, we must multiply by its inverse instead. But inverting a matrix can be expensive!

So we can move it to the other side of the equation and solve this linear system instead. However, the Hessian itself can be quite expensive to compute!

# Why Gauss-Newton?

Gradient Descent?

$$\mathbf{a}_{n+1} = \mathbf{a}_n$$



Newton's Method

$$x_{n+1} = x_n - \frac{f'(x_n)}{f''(x_n)}$$

Single variable...

$$\mathbf{x}_{n+1} = \mathbf{x}_n - [\mathbf{H}f(\mathbf{x}_n)]^{-1} \nabla f(\mathbf{x}_n)$$

$$[\mathbf{H}f(\mathbf{x}_n)] \Delta \mathbf{x} = -\nabla f(\mathbf{x}_n)$$

Why are we doing Gauss-Newton? (and what is it anyway?)

Why don't we use gradient descent for optimization like all the cool kids? It's simple to understand, and obviously works.

Well, it actually has bad convergence on interesting functions. There is a reason even beyond sheer data size why neural nets take forever to train.

But simplicity is good! The problem with gradient descent is we aren't using much information, we are linearly approximating our function at each step to find our next move. If we quadratically approximate instead, we get Newton's method for optimization, which is taught as a root finding method in high schools all over. That seems promisingly simple, while having better convergence properties. Of course what I put up there has a problem, its one dimensional.

We can move to higher dimensions through analogy. The derivative becomes the gradient vector, and the second derivative is the Hessian matrix.

Problem, we can't divide by a matrix, we must multiply by its inverse instead. But inverting a matrix can be expensive!

So we can move it to the other side of the equation and solve this linear system instead. However, the Hessian itself can be quite expensive to compute!

## Gauss Newton as Approximate Newton

$$[\mathbf{H}f(\mathbf{x}_n)]\Delta\mathbf{x} = -\nabla f(\mathbf{x}_n)$$

$$f(\mathbf{x}) = \sum_{i=1}^m r_i^2(\mathbf{x})$$

$$H_{jk} = 2 \sum_{i=1}^m \left( \frac{\partial r_i}{\partial x_j} \frac{\partial r_i}{\partial x_k} + r_i \frac{\partial^2 r_i}{\partial x_j \partial x_k} \right)$$

$$H_{jk} \approx 2 \sum_{i=1}^m \frac{\partial r_i}{\partial x_j} \frac{\partial r_i}{\partial x_k} \quad \longrightarrow \quad \mathbf{J}^T \mathbf{J} \Delta\mathbf{x} = -\mathbf{J}^T \mathbf{f}$$

$$H_{jk} \approx 2 \sum_{i=1}^m J_{ij} J_{ik}$$

We have Newton's method.

In the special case of the function  $f$  being a sum of squared residual terms, we can exploit the extra structure. If we write out an element of the Hessian matrix we get this sum of terms dependent on a single residual (and its derivatives).

The trick is if we drop the higher order terms on the right, we are simply using the Jacobian transpose Jacobian, which is often a much easier to compute value.

## One Step of Gauss-Newton

$$\mathbf{J}^T \mathbf{J} \Delta \mathbf{x} = -\mathbf{J}^T \mathbf{f}$$

$\Delta \mathbf{x}$  = Step to apply to unknowns

$\mathbf{r}$  = vector of residuals

$\mathbf{J}$  = Jacobian of  $\mathbf{r}$  with respect to unknowns

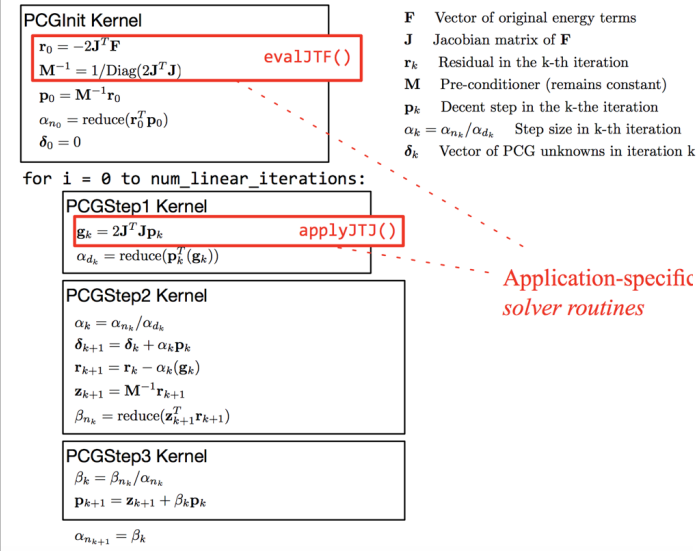


## Generic Gauss Newton

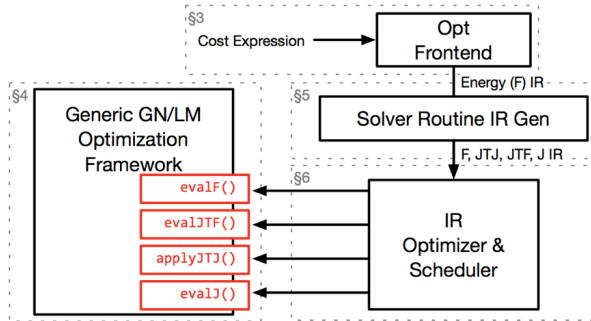
while (nonlinear convergence criteria false):

    Solve  $\mathbf{J}^T \mathbf{J} \Delta \mathbf{x} = -\mathbf{J}^T \mathbf{r}$  for  $\Delta \mathbf{x}$

$\mathbf{x} = \mathbf{x} + \Delta \mathbf{x}$



Application-specific solver routines



## Laplacian Smoothing bandwidth SoL

GeForce 980 ~224GB/s of bandwidth, 1MP image

- Full matrix: ~20.1s touch every element
- CRS matrix: ~0.17ms touch every element
- Matrix-free: ~0.017ms touch every pixel

(d) Representation of non-zero entries in the expression  $\mathbf{g} = 2\mathbf{J}^T\mathbf{J}\mathbf{p}$  that are required to calculate  $g_{0,0}$

$$\begin{array}{c} \text{unknowns} \uparrow \\ \mathbf{g} \\ \left[ \begin{array}{c} g_{0,0} \end{array} \right] \end{array} \quad \begin{array}{c} \leftarrow \text{required row} \rightarrow \\ \\ \\ \\ \\ \\ \\ \\ \text{unknowns} \uparrow \\ \text{residuals} \rightarrow \end{array} \quad = \quad 2 \quad \begin{array}{c} \text{residuals} \rightarrow \\ \left[ \begin{array}{ccccc} \text{fit}_{0,0} & \text{h\_reg}_{0,0} & \text{h\_reg}_{-1,0} & \text{v\_reg}_{0,0} & \text{v\_reg}_{0,-1} \\ \frac{d\text{fit}_{0,0}}{dx_{0,0}} & \frac{dh\_reg_{0,0}}{dx_{0,0}} & \frac{dh\_reg_{-1,0}}{dx_{0,0}} & \frac{dv\_reg_{0,0}}{dx_{0,0}} & \frac{dv\_reg_{0,-1}}{dx_{0,0}} \\ \\ \text{Row corresponding to } g_{0,0} \text{ has non-zeros} \\ \text{for each residual containing } x_{0,0} \\ \\ \\ \\ \text{Rows are required for each} \\ \text{non-zero column required in } \mathbf{J}^T \end{array} \right] \end{array}$$

$$\begin{array}{c} \text{residuals} \rightarrow \\ \left[ \begin{array}{ccccc} \text{fit}_{0,0} & \text{h\_reg}_{0,0} & \text{h\_reg}_{-1,0} & \text{v\_reg}_{0,0} & \text{v\_reg}_{0,-1} \\ \frac{d\text{fit}_{0,0}}{dx_{0,0}} & \frac{dh\_reg_{0,0}}{dx_{0,0}} & \frac{dh\_reg_{-1,0}}{dx_{0,0}} & \frac{dv\_reg_{0,0}}{dx_{0,0}} & \frac{dv\_reg_{0,-1}}{dx_{0,0}} \\ \frac{dh\_reg_{0,0}}{dx_{1,0}} & \frac{dh\_reg_{0,0}}{dx_{-1,0}} & & & \\ \frac{dh\_reg_{-1,0}}{dx_{0,0}} & & \frac{dh\_reg_{-1,0}}{dx_{-1,0}} & & \\ \frac{dv\_reg_{0,0}}{dx_{0,0}} & & & \frac{dv\_reg_{0,0}}{dx_{0,1}} & \\ \frac{dv\_reg_{0,-1}}{dx_{0,0}} & & & & \frac{dv\_reg_{0,-1}}{dx_{0,-1}} \end{array} \right] \end{array} \quad \begin{array}{c} \text{unknowns} \rightarrow \\ \mathbf{J} \\ \left[ \begin{array}{ccccc} x_{0,0} & x_{1,0} & x_{-1,0} & x_{0,1} & x_{0,-1} \end{array} \right] \end{array}$$

$$\begin{array}{c} \text{unknowns} \uparrow \\ \mathbf{p} \\ \left[ \begin{array}{c} p_{0,0} \\ p_{1,0} \\ p_{-1,0} \\ p_{0,1} \\ p_{0,-1} \end{array} \right] \end{array} \quad \begin{array}{c} \text{unknowns} \rightarrow \\ \left[ \begin{array}{c} x_{0,0} \\ x_{1,0} \\ x_{-1,0} \\ x_{0,1} \\ x_{0,-1} \end{array} \right] \end{array}$$

non-zero columns where each individual residual has support

## Miscellaneous

- Levenberg-Marquardt
- Cuspars backend
- Float/Double precision
- IRLS – Can solve for  $L_1$  energies
- OpenGL-like C-API