



# AD Master Class 4

Pushing performance using  
SIMD vectorization

**Viktor Mosenkis**  
`viktor.mosenkis@nag.co.uk`



Experts in numerical algorithms  
and HPC services

*20 August 2020*

# AD Masterclass Schedule and Remarks

---

## ■ AD Masterclass Schedule

- ☐ 30 July 2020 | Why the need for Algorithmic Differentiation?
- ☐ 6 August 2020 | How AD works
- ☐ 13 August 2020 | Testing and validation
- ☐ 20 August 2020 | Pushing performance using SIMD vectorization
- ☐ 27 August 2020 | Bootstrapping validated adjoints on real-world codes

## ■ Remarks

- ☐ Please submit your questions via the questions panel at any time during this session, these will be addressed at the end.
- ☐ A recording of this session, along with the slides will be shared with you in a day or two.

# Dialogue

---

We want this webinar series to be interactive (even though it's hard to do)

- We want your feedback, we want to adapt material to your feedback
- Please feel free to contact us via email to ask questions at any time
- We'd love to reach out offline, discuss what's working, what to spend more time on
- For some orgs, may make sense for us to do a few bespoke sessions

Blog:

[https://www.nag.com/blog/  
algorithmic-differentiation-masterclass-series-2](https://www.nag.com/blog/algorithmic-differentiation-masterclass-series-2)

# Outcomes

---

- Understand the idea of the vector versions for both AD models
  - tangent-linear and
  - adjoint model
- Learn how to compute the Jacobian using tangent-linear and adjoint vector model with `dco/c++`
- Learn how vector models in conjunction with AVX/SIMD instruction can speed up the computation of the Jacobian
- Continue discussion on which AD model should be used

# Algorithmic Differentiation: Tangent-Linear Model

---

$$F : \mathbb{R}^n \rightarrow \mathbb{R}^m, \quad \mathbf{y} = F(\mathbf{x})$$

## ■ Tangent-Linear (**scalar**) Model (TLM) $\dot{F}$

$$\dot{F}(\mathbf{x}, \dot{\mathbf{x}}) = \underset{\in \mathbb{R}^{m \times n}}{F'(\mathbf{x})} \cdot \underset{\in \mathbb{R}^n}{\dot{\mathbf{x}}} = \underset{\in \mathbb{R}^m}{\dot{\mathbf{y}}}$$

$$\square \quad \frac{Cost(\dot{F})}{Cost(F)} \approx 2$$

## ■ Tangent-Linear (**vector**) Model (TLM) $\dot{F}_k$

$$\dot{F}_k(\mathbf{x}, \dot{\mathbf{X}}) = \underset{\in \mathbb{R}^{m \times n}}{F'(\mathbf{x})} \cdot \underset{\in \mathbb{R}^{n \times k}}{\dot{\mathbf{X}}} = (F'(\mathbf{x}) \cdot \dot{\mathbf{x}}_1, \dots, F'(\mathbf{x}) \cdot \dot{\mathbf{x}}_k) = \underset{\in \mathbb{R}^{m \times k}}{\dot{\mathbf{Y}}}$$

$$\square \quad \frac{Cost(\dot{F}_k)}{Cost(F)} \approx k + 1$$

# Jacobian with tangent-linear vector model

---

- In tangent linear **scalar** model we compute the **Jacobian column-wise** by setting  $\dot{\mathbf{x}}$  to Cartesian basis vectors of  $\mathbb{R}^n$   $e_i : 1 \leq i \leq n$ . Especially

$$\dot{F}(\mathbf{x}, e_i) = F'(\mathbf{x}) \cdot e_i = F'(\mathbf{x})_{-,i},$$

where  $F'(\mathbf{x})_{-,i}$  denotes the  $i$ -th column of  $F'(\mathbf{x})$ .

- In tangent linear **vector** model we compute  $k$  **columns** of the **Jacobian simultaneously** by setting  $\dot{\mathbf{X}} = I_n^{i,k}$ , where  $I_n^{i,k} = (e_i, e_{i+1}, \dots, e_{i+k})$  ( an  $n \times k$  block identity matrix). Hence

$$\dot{F}_k(\mathbf{x}, I_n^{i,k}) = F'(\mathbf{x}) \cdot I_n^{i,k} = F'(\mathbf{x})(e_i, \dots, e_{i+k}) = F'(\mathbf{x})_{-,i\dots i+k}$$

where  $F'(\mathbf{x})_{-,i\dots i+k}$  denotes the  $i$ -th to  $i+k$ -th column of  $F'(\mathbf{x})$ .

# Inside tangent-linear vector model

$$F : \mathbb{R}^2 \rightarrow \mathbb{R}, \quad y = \sin(x_1 \cdot x_2)$$

**Tangent linear code** (2 directions)

$$\dot{u}^1 = x_2 \cdot \dot{x}_1^1 + x_1 \cdot \dot{x}_2^1$$

$$\dot{u}^2 = x_2 \cdot \dot{x}_1^2 + x_1 \cdot \dot{x}_2^2$$

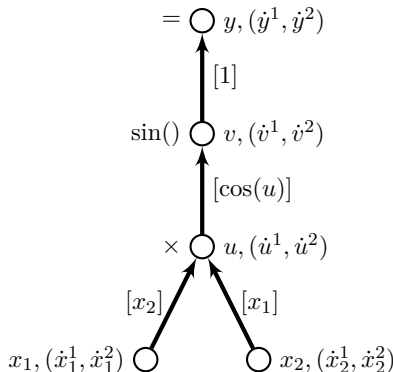
$$u = x_1 \cdot x_2$$

$$\dot{v}^1 = \cos(u) \cdot \dot{u}^1 \quad \dot{v}^2 = \cos(u) \cdot \dot{u}^2$$

$$v = \sin(u)$$

$$\dot{y}^1 = \dot{v}^1 \quad \dot{y}^2 = \dot{v}^2$$

$$y = v$$



# Algorithmic Differentiation: Adjoint Model

---

$$F : \mathbb{R}^n \rightarrow \mathbb{R}^m, \quad \mathbf{y} = F(\mathbf{x})$$

## ■ Adjoint Model (**scalar**) (ADM) $\bar{F}$ (reverse mode)

$$\bar{F}(\mathbf{x}, \bar{\mathbf{y}}) = \underbrace{\bar{\mathbf{y}}}_{\in \mathbb{R}^m} \cdot \underbrace{F'(\mathbf{x})}_{\in \mathbb{R}^{m \times n}} = F'(\mathbf{x})^T \cdot \bar{\mathbf{y}} = \bar{\mathbf{x}}$$

$$\square \quad \frac{\text{Cost}(\bar{F})}{\text{Cost}(F)} = M$$

## ■ Adjoint (**vector**) Model (ADM) $\bar{F}_k$

$$\bar{F}_k(\mathbf{x}, \bar{\mathbf{Y}}) = F'(\mathbf{x})^T \cdot \underbrace{\bar{\mathbf{Y}}}_{\in \mathbb{R}^{m \times k}} = (F'(\mathbf{x})^T \cdot \bar{\mathbf{y}}_1, \dots, F'(\mathbf{x})^T \cdot \bar{\mathbf{y}}_k) = \bar{\mathbf{X}}$$

$$\square \quad \frac{\text{Cost}(\bar{F}_k)}{\text{Cost}(F)} = O(k) < M \cdot (k)$$



# Jacobian with adjoint vector model

---

- In adjoint **scalar** model we compute the **Jacobian row-wise** by setting  $\bar{\mathbf{y}}$  to Cartesian basis vectors of  $\mathbb{R}^m$   $e_i : 1 \leq i \leq m$ . Especially

$$\bar{F}(\mathbf{x}, e_i) = F'(x)^T \cdot e_i = F'(x)_{i,-},$$

where  $F'(x)_{i,-}$  denotes the  $i$ -th row of  $F'(x)$ .

- In adjoint **vector** model we compute  $k$  rows of the **Jacobian simultaneously** by setting  $\bar{\mathbf{Y}} = I_m^{i,k}$ , where  $I_m^{i,k} = (e_i, e_{i+1}, \dots, e_{i+k})$  ( an  $m \times k$  block identity matrix). Hence

$$\bar{F}_k(\mathbf{x}, I_n^{i,k}) = F'(x)^T \cdot I_n^{i,k} = F'(x)^T (e_i, \dots, e_{i+k}) = F'(x)_{i\dots i+k,-}$$

where  $F'(x)_{i\dots i+k,-}$  denotes the  $i$ -th to  $i+k$ -th row of  $F'(x)$ .

# Inside adjoint vector model

$$F : \mathbb{R}^2 \rightarrow \mathbb{R}, \quad y = \sin(x_1 \cdot x_2)$$

## Adjoint Code (2 directions)

Forward Sweep (primal computation)

$$u = x_1 \cdot x_2$$

$$v = \sin(u)$$

$$y = v$$

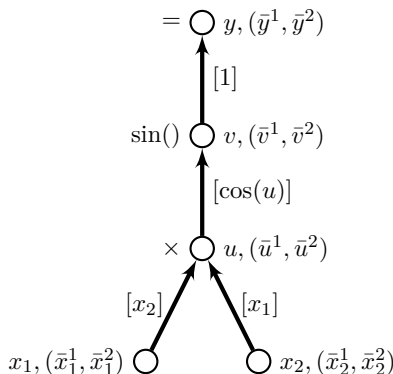
Reverse Sweep (adjoint computation)

$$\bar{v}^1 = \bar{y}^1 \quad \bar{v}^2 = \bar{y}^2$$

$$\bar{u}^1 = \cos(u) \cdot \bar{v}^1 \quad \bar{u}^2 = \cos(u) \cdot \bar{v}^2$$

$$\bar{x}_1^1 = x_2 \cdot \bar{u}^1 \quad \bar{x}_1^2 = x_2 \cdot \bar{u}^2$$

$$\bar{x}_2^1 = x_1 \cdot \bar{u}^1 \quad \bar{x}_2^2 = x_1 \cdot \bar{u}^2$$



# Tangent-linear vector Model with dco/c++

---

- Replace floating point variables with `dco::gt1v<double, k>::type`, where `k` is the size of the tangent vector
- Write the driver
- Conceptually `dco::gt1v<double, k>::type` contains two components, both components are computed during normal execution
  - ☐ value
  - ☐ tangent (array)
- Interface
  - ☐ `dco::value(DCO_TYPE)` - access to value component
  - ☐ `dco::derivative(DCO_TYPE)` - access to the tangent (array) component

## Tangent-Linear vector Model: Jacobian with dco/c++

---

```
1  template <typename T>
2  void foo(int& n, T* x, T& y){ ... }
3  int main(){
4      DCO_TANGENT_VECTOR_TYPE *x, y;
5      ...
6      for (int i = 0; i < n; i+=k) {
7          for (int j = 0; j < k; j++)
8              dco::derivative(x[i+j])[j] = 1.0;
9
10         foo(n, x, y);
11
12         for (int j = 0; j < k; j++)
13             J[i+j] = dco::derivative(y)[j];
14         for (int j = 0; j < k; j++)
15             dco::derivative(x[i+j])[j] = 0.0;
16     }
```

The tangent-linear vector model of `foo` is executed  $n/k$  times

# Tangent-linear vector Model

---

## Why should I use it?

- each time tangent-linear (scalar/vector) model is called the function value is computed along with tangent-linear projections.
- with tangent-linear vector model we perform only  $n/k$  primal function evaluations instead of  $n$  with tangent-linear scalar model
- computing several projections at the same time can profit from **AVX/SIMD** extensions on your CPU.

## Problems

- higher memory requirements to store the additional tangents
- more complicated driver

# Tangent-Linear vector Model

---

**How should I choose the vector size for  $gt1v$ ?**

- Choosing the vector size as big as possible ( $k=n$ )
  - ☐ reduces the number of primal function evaluations
  - ☐ high memory requirements. Each  $gt1v$  variable contains  $n + 1$  doubles
  - ☐ bad memory access pattern (caching effects)
  
- On modern Intel CPU's with AVX2/AVX512 vector sizes (8/12/16) is almost always the best choice.
  - ☐ choose vector sizes as multiple of four/eight to be aligned with AVX2/AVX512
  - ☐ the unused directions should be set to zero (number of inputs not multiple of the chosen vector size).

## Computing several adjoint directions with `ga1s`

---

```
1 void foo(const int &n, const int &m, T* x, T* y) {...}
2 int main(){
3     DCO_ADJOINT_VECTOR_TYPE *x, *y;
4     ...
5     foo(n, m, x, y); //Record the tape
6     for (int i = 0; i < m; i++) { //m tape interpretations
7         dco::derivative(y[i]) = 1.0; //set adjoint of y
8
9         DCO_M::global_tape->interpret_adjoint();
10
11        for (int j = 0; j < n; j++) //Extract the derivatives
12            J[i*n + j] = dco::derivative(x[j]);
13        //zero all adjoints
14        DCO_M::global_tape->zero_adjoints();
15    }
16 }
```

Function `foo` is recorded once but is interpreted *m* times

## Computing several adjoint directions: Remarks

---

- Record the function only once and interpret as many times as you need
- Interpretation step is much cheaper compared to recording step. Only 10%-30% of the overall time.
  - during recording the DAG (tape) is written to memory including computation of partial derivatives
  - tape interpretation performs only fma's operations. DAG (tape) is left untouched only the adjoint vector is updated.
- Vector adjoint model - can save only the interpretation time. Recording time is the same.

NOTE: All information only applies to adjoint model not using advanced adjoint techniques (e.g. checkpointing, symbolic adjoints)



# Adjoint vector Model with dco/c++

---

- Replace floating point variables with `dco::ga1v<double, k>::type`, where `k` is the number of the adjoint directions computed simultaneously
- Write the driver
- Conceptually `dco::ga1v<double, k>::type` contains two components
  - ☐ value
  - ☐ adjoint (array)

During the execution of the function `dco/c++` computes the value component and records the computational graph (tape). Interpretation of the tape is needed to compute the adjoint components.

# Adjoint vector Model with dco/c++: Basic Interface

---

## ■ Interface of `dco::galv<double, k>::type`

- `dco::value(DCO_TYPE)` - access to value component
- `dco::derivative(DCO_TYPE)` - access to the adjoint (array) component

## ■ Interface of the tape `DCO_MODE::tape_t`

- `DCO_TAPE_TYPE::create()` - creates tape and returns pointer to it
- `DCO_TAPE_TYPE::remove(DCO_TAPE_TYPE*)` - deallocates tape
- `DCO_TAPE_TYPE::register_variable(DCO_TYPE)` - Marks variable as independent
- `DCO_TAPE_TYPE::register_output_variable(DCO_TYPE)` - Marks variable as dependent
- `DCO_TAPE_TYPE::interpret_adjoint()` - Runs tape interpreter

## Adjoint vector Model with dco/c++: Jacobian

---

```
1  template <typename T>
2  void foo(int& n, int& m, T* x, T& y){ ... }
3  int main(){
4      DCO_ADJOINT_VECTOR_TYPE *x, y;
5      ...
6      for (int i = 0; i < m; i+=k) {
7          for (int j = 0; j < k; j++)
8              dco::derivative(y[i+j])[j] = 1.0;
9          DCO_M::global_tape->interpret_adjoint(); //Interpret
10         for (int s = 0; s < k; s++) {
11             for (int j = 0; j < n; j++)
12                 J[(i+s) * n + j] = dco::derivative(x[j])[s];
13         }
14         DCO_M::global_tape->zero_adjoints();
15     }
16 }
```

Function `foo` is recorded once but is interpreted  $m/k$  times

# Adjoint vector Model

---

## Why should I use it?

- with adjoint vector model we perform only  $m/k$  tape interpretations instead of  $m$  with adjoint scalar model. Can improve performance when used in conjunction with advanced adjoint techniques (checkpointing, symbolic adjoints)
- computing several adjoint projections at the same time can profit from **AVX/SIMD** extensions on your CPU.

## Problems

- higher memory requirements to store the additional adjoints directions. Can be addressed by using **compressed adjoint vector** (Reuse adjoint storage by analysing the maximum number of required distinct adjoint memory locations (Master Class 5)).
- more complicated driver

# Adjoint vector Model

---

**How should I choose the vector size for `ga1v`?**

- Choosing the vector size as big as possible ( $k=m$ )
  - ☐ reduces the number of tape interpretations
  - ☐ high memory requirements (Each entry in the vector of adjoints is now a vector of doubles with size  $k$ .)
  - ☐ bad memory access pattern (caching effects)
- On modern Intel CPU's with AVX2/AVX512 vector sizes (4/8/12/16) is almost always the best choice.
  - ☐ choose vector sizes as multiple of four/eight to be aligned with AVX2/AVX512
  - ☐ the unused directions should be set to zero (number of outputs not multiple of the chosen vector size).

# Tangents vs Adjoint

---

In Masterclass 2 we already learned that if  $M \cdot m < 2 \cdot n$  you should use the adjoint model and else the tangent-linear model. Now let's see how the tangent vector model can change this statement:

## ■ Pro tangent

- Tangent vector mode can speed up your tangent computation by a factor of  $\approx 3$  on a AVX2 machine. AVX512 should provide even better results.
- Hence for  $m = 1$  we should rather say that  $M < n$  to justify the usage of adjoints
- For  $m > 1$  please note that adjoint model can be very efficient on computing several rows of Jacobian (see next slide).

# Tangents vs Adjoints

---

In Masterclass 2 we already learned that if  $M \cdot m < 2 \cdot n$  you should use the adjoint model and else the tangent-linear model. Now let's see how the adjoint vector model can change this statement:

## ■ Pro adjoint

- Computing additional directions with adjoint **scalar** model is  $\approx 10\% - 30\%$  overall runtime. (**without advanced adjoint techniques**)
- Using adjoint vector mode provides an addition speed up of  $\approx 3$  on a AVX2 machine compared to the adjoint scalar model. AVX512 machine should perform even better.
- Even for big  $m$  adjoint might be the right choice, if the tape is **small enough**. Especially if  $M < 10$ .

# Tangents vs Adjoint

---

- Consider using all available models and know their strength and weaknesses. Don't do adjoint because somebody tells you it is great.
- Start with the implementation of the tangent model (you need it anyway as discussed in Master Class 3)
- Based on performance on tangent-linear model you can decide, whether you need the adjoint model
  - $M > 30$  - naive implementation of adjoints should give you desired performance
  - $M < 20$  - you will need some tweaks (e.g. symbolic adjoints)
  - $M < 10$  - very good adjoint code significant development time required
- Adjoint requires control flow reversal, making it challenging for parallel codes and increases the memory usage (forcing to use advanced adjoint techniques).



# Summary

---

In this Class we learned

- How tangent-linear and adjoint vector model work.
- How they can be used to speed up the computation
- How AVX/SIMD extension can help us speed up
  - ☐ tangent-linear vector Model
  - ☐ adjoint vector Model
- Continue discussion on tangent vs. adjoints

### **Bootstrapping validated adjoints on real-world codes**

In the next part our we will learn different automatic techniques to avoid running out of memory when computing the adjoints such as

- write tape to disk
- Jacobian pre-accumulation
- compressing the vector of adjoints

The goal to compute an adjoint result that we can use as a reference for potential later code optimisations.

You will see a survey on your screen after exiting  
from this session.

We would appreciate your feedback.

We are now moving on the Q&A Session