

# Large-scale semantic classification:

## Outcome of the first year of Inria aerial image labeling benchmark

B. Huang<sup>1</sup>, K. Lu<sup>2</sup>, N. Audebert<sup>3,4</sup>, A. Khalef<sup>5</sup>, Y. Tarabalka<sup>6</sup>, J. Malof<sup>1</sup>, A. Boulch<sup>3</sup>, B. Le Saux<sup>3</sup>, L. Collins<sup>1</sup>, K. Bradbury<sup>1</sup>, S. Lefèvre<sup>4</sup>, M. El-Saban<sup>5</sup>

<sup>1</sup> Duke University; <sup>2</sup> NUS; <sup>3</sup> ONERA; <sup>4</sup> Univ. Bretagne-Sud, IRISA; <sup>5</sup> Raisa energy; <sup>6</sup> UCA, Inria.

⇒ [project.inria.fr/aerialimagelabeling/](http://project.inria.fr/aerialimagelabeling/)

### Inria Benchmark dataset and statistics

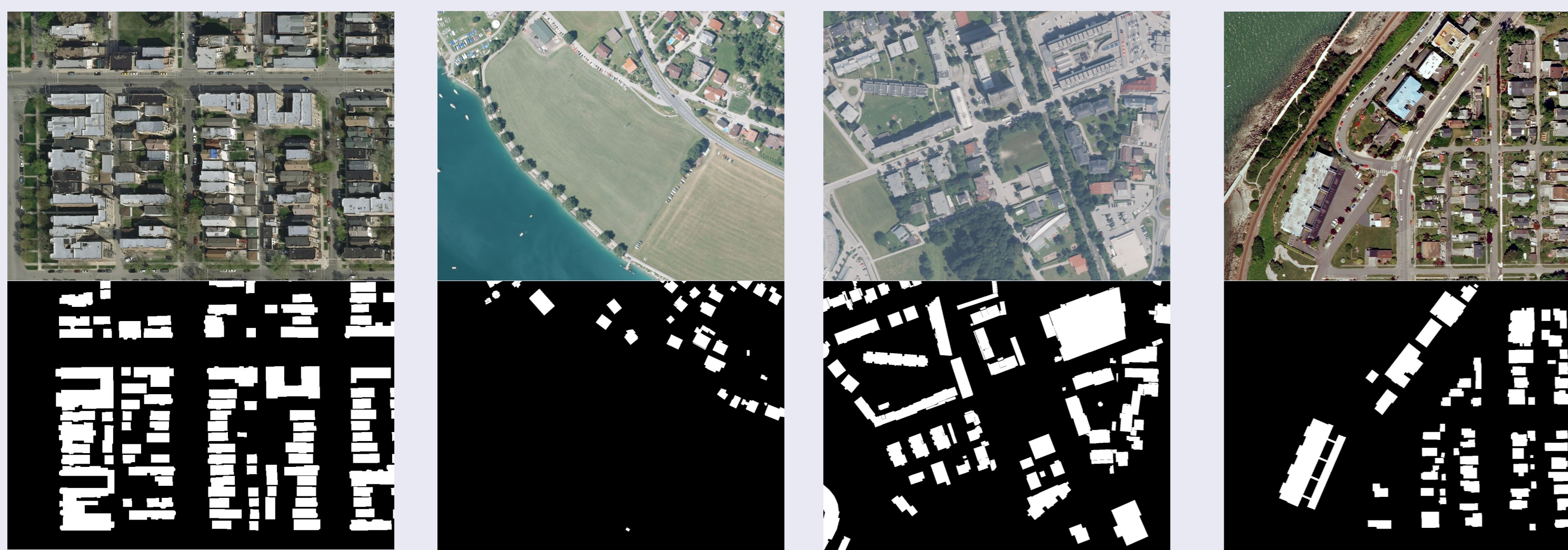
- Problem: Large-scale pixelwise semantic labeling of aerial images
- Two semantic classes: *building* and *not building* (ref. data by rasterizing building footprints)
- Different cities in train and test subsets  
⇒ E.g., we should classify San Francisco without “seeing” it before
- European/American & high-/low-density urban landscapes in both subsets
- 0.3 m spatial resolution, 3 color bands, 360 tiles (1500<sup>2</sup> px each)

Statistics:

Train	Tiles	Total area	Test	Tiles	Total area
Austin, TX	36	81 km <sup>2</sup>	Bellingham, WA	36	81 km <sup>2</sup>
Chicago, IL	36	81 km <sup>2</sup>	San Francisco, CA	36	81 km <sup>2</sup>
Kitsap County, WA	36	81 km <sup>2</sup>	Bloomington, IN	36	81 km <sup>2</sup>
Vienna, Austria	36	81 km <sup>2</sup>	Innsbruck, Austria	36	81 km <sup>2</sup>
West Tyrol, Austria	36	81 km <sup>2</sup>	East Tyrol, Austria	36	81 km <sup>2</sup>
<b>Total</b>	<b>180</b>	<b>405 km<sup>2</sup></b>	<b>Total</b>	<b>180</b>	<b>405 km<sup>2</sup></b>

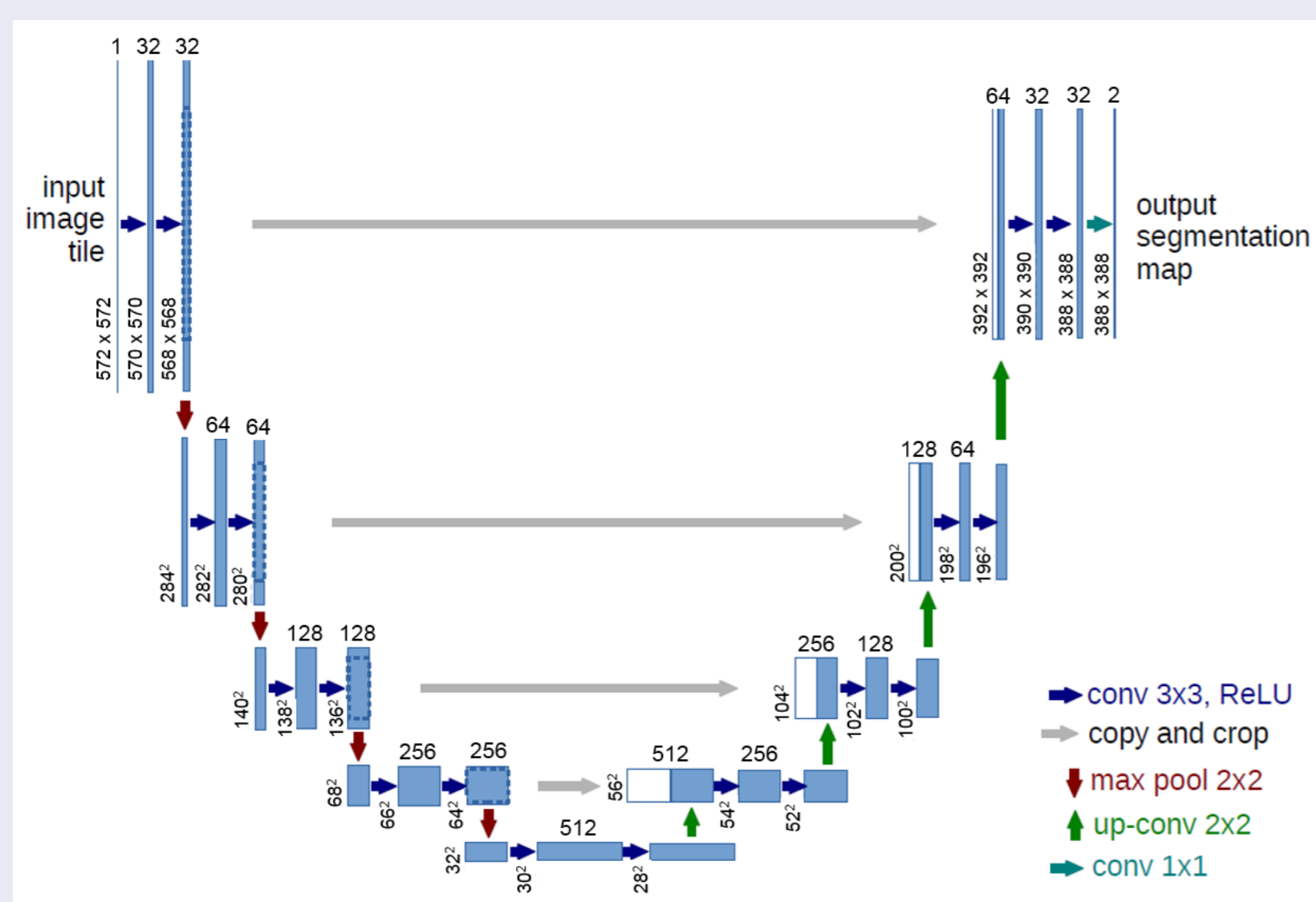
- During the first year after the benchmark release:
  - > 800 downloads from all continents, from public & private institutes
  - 16 submissions with the results on the test set
  - Which method is the best?** Four winning methods are detailed here

### Close-ups of training and test sets



Chicago (train) West Tyrol (train) East Tyrol (test) Bellingham (test)

### 1st place: U-net with novel training/test strategy (AMLL, Duke Univ.)



- Original **U-net** architecture [1] with half as many filters at each layer
- Training strategy:**
  - From training dataset: tiles 6-36 from each city for training, the rest for validation
  - Extract 572 × 572 input patches on a uniform grid, with 92 pixels of overlap between neighboring patches
  - Minibatch of 5 randomly selected patches
  - Data augmentation: vertical/horizontal flips and orthogonal rotations
  - Cross-entropy objective function
  - Adam optimizer: initial learning rate of 1e − 3, a momentum of 0.9
  - 100 epochs, each epoch processes 8000 minibatches
- Label inference:**
  - U-net predicts poorly at the edge of its output
  - To mitigate this problem ⇒ use 2636 × 2636 input patches\* during label inference

[1] Ronneberger et al., “U-net: Convolutional networks for biomedical image segmentation,” in *MICCAI*, 2015.

\* Maximum size supported by 1080 Ti GPU.

### 2nd place: Dual-resolution U-net (NUS)

- U-net** architecture with a pair of dual-resolution images as input
  - Crop high-resolution 384 × 384 patches
  - Crop 768 × 768 patches with the same center and downsample them to 384 × 384 patches
  - Features from high & low resolution patches are extracted by U-net
  - Result = weighted sum of dual-resolution score maps
- Loss function** = combination of sigmoid cross-entropy (*sigmCE*) and a soft Jaccard loss [2]:

$$L_{NUS} = L_{sigmCE} - \log I_{soft-IJU}$$

- Implementation details:**
  - Channels of the modified U-net are: 32, 64, 128, 128, 256, 128, 128, 64, 32
  - Data augmentation: vertical/horizontal flips
  - Adam optimizer: initial learning rate of 1e − 3, a momentum of 0.9, “poly” learning rate policy
  - 30 epochs

[2] Mattyus et al., “Deeproadmapper: Extracting road topology from aerial images,” in *ICCV*, 2017.

### 3rd place: Signed distance transform regression (ONERA)

- Standard **SegNet** architecture with pre-trained VGG-16 weights
  - 384 × 384 patches, stochastic gradient descent optimizer
- Include spatial context in optimization  
⇒ Add a regularization loss computed on the Euclidean signed distance transform (SDT) [3]:

$$L_{ONERA} = NLLLoss(Z_{seg}, Y_{seg}) + \lambda L1(Z_{dist}, Y_{dist}),$$

where *NLLLoss* = negative log-likelihood loss function, *L1* = L1 penalty on SDT distances,  $\lambda$  = hyper-parameter

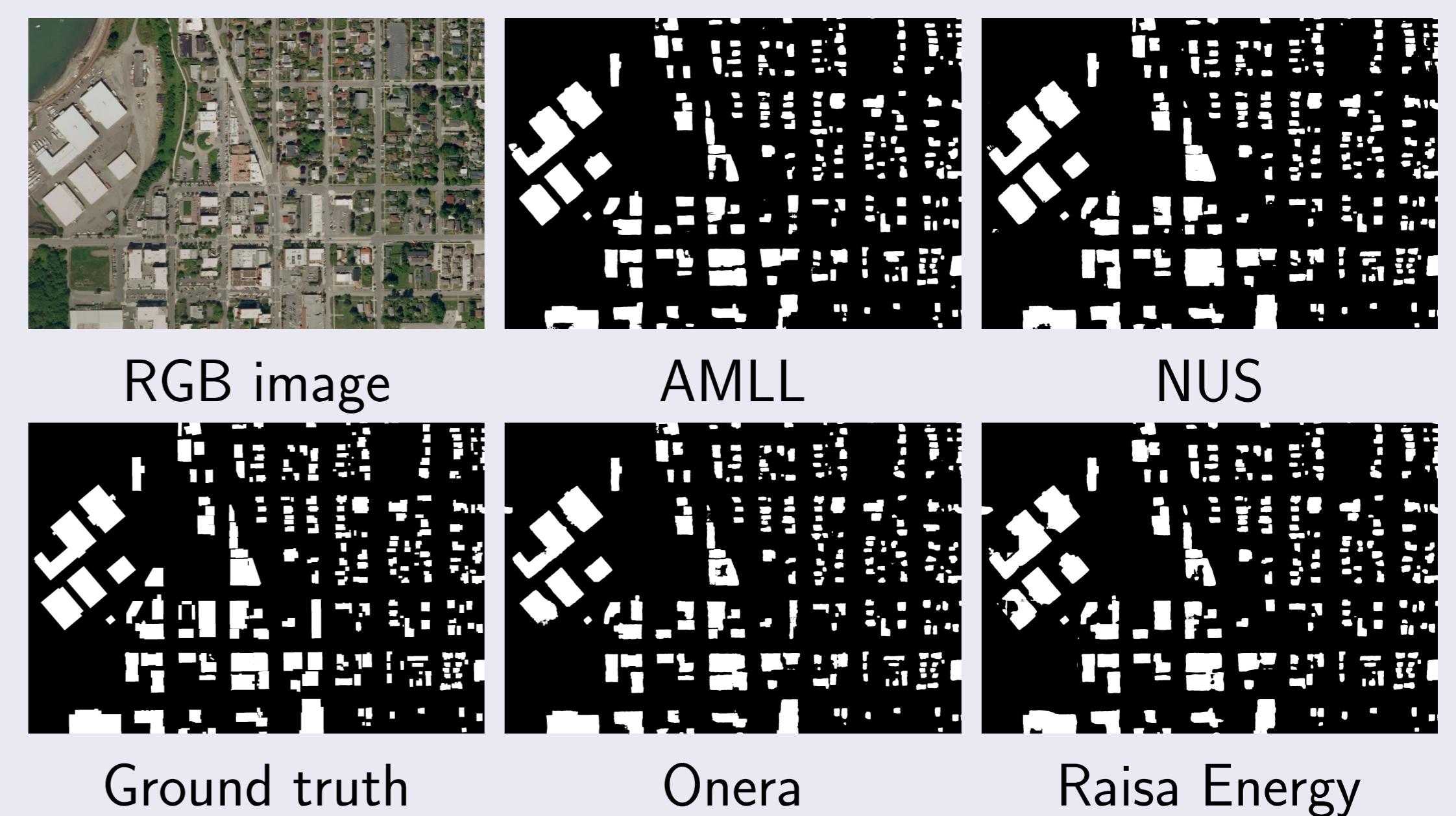
[3] Ye, “The signed Euclidean distance transform and its applications,” in *ICPR*, 1988.

### 4th place: Stacked U-nets (Raisa Energy)

- Stack of two **U-nets** arranged end-to-end
  - Second net enhances predictions of the first net
- Loss function** combines binary cross entropy and a differential form of Intersection over union (IoU) [2]

### Experimental results

- <https://project.inria.fr/aerialimagelabeling/leaderboard/>



	Belling.	Bloom.	Inns.	S. Francisco	East Tyrol	Overall
AMLL	67.14	65.43	72.27	<b>75.72</b>	74.67	<b>72.55</b>
NUS	<b>70.74</b>	66.06	<b>73.17</b>	73.57	<b>76.06</b>	72.45
ONERA	68.92	<b>68.12</b>	71.87	71.17	74.75	71.02
RAISA	68.73	60.83	70.07	70.64	74.76	69.57

Numerical evaluation on test set (IoU scores)

### Concluding remarks

- Active exploitation of the benchmark since its release
- U-net architecture has shown the highest performance
- Good choice of loss function & training strategy boosts results
- Published on Nov '16, > 1500 downloads as of June '18, > 50 submissions to contest
- Contest still open** to submit results to benchmark!