

# WormSeg User Guide

Kyle Moy

February 2, 2016

## 1 Introduction

This iteration of the worm segmentation software was rewritten from scratch and includes improvements to usability, scalability, and performance.

The software comprises a single Java package that can be configured via command line arguments. This makes it possible to launch and automate the software via scripts, and will require no human input once launched. The software will manage its own processes, thus only one instance needs to be running on any individual machine or VM.

This iteration is built on a different data-scattering model that, at the cost of memory, completely eliminates disk I/O on nodes. This inherently places a hard limitation on the software, requiring the nodes to be able to fit the its entire workload in memory<sup>1</sup>. Disk I/O was responsible for ~40% of the total runtime, and its elimination reflects significant improvement. A secondary characteristic of this model is that many more nodes can be effectively utilized, increasing the scalability of this algorithm. Under the previous model, there is a point where transferring data to more nodes would cost more than the computation time (although this is true of all parallel algorithms).

Finally, since I was rewriting from scratch, I had the opportunity to make major changes to the algorithm as well as include micro-optimizations from the beginning. With the exception of substituting box blur for Gaussian blur, no image processing techniques have been changed. Optimizations mostly involve the way data is represented internally and throwing away a lot of unused data (ARGB data reduced to only the blue channel, chosen arbitrarily). Micro-optimizations turned out to be significant, for example pre-calculating division tables, looping arrays in memory-sequential order, memoizing certain lookups, etc. which shaved off several milliseconds per frame. Although I don't have many trial cases, my tests show about 70% reduction in time-per-frame per thread compared to the previous software.

---

<sup>1</sup>Our largest input data to date, egl19-f2, occupies 5.30GB, which fits in MEDIXSRV memory. Unlike the host, nodes only hold their fraction of the whole payload.

## 2 Instructions

1. Configure a launch script for both the host and node. Refer to the Launch Arguments section.
2. Launch the host and any number of nodes, in any order. Host will begin when enough nodes have joined. (See host argument "-nodes")
3. Wait for the task to complete. Host will output the results and then terminate.
4. Nodes will NOT terminate on their own. This allows multiple tasks to be run by launching hosts repeatedly (but not concurrently).

## 3 Launch Arguments

Bracketed arguments have default values and thus are optional. Strings and numbers are entered literally and should not be enclosed in quotes.

### 3.1 Host Arguments

#### 3.1.1 Operational Arguments

**-mode (host / node)**

This must be specified with "**node**" to start the program in host mode.

**-project (name)**

The name of the project to process. This argument, along with "-dir", will be used to detect the input and output paths. This argument **must** be specified if "-input" and "-output" are not.

**[-dir (path)]**

The path<sup>2</sup> to the directory where project folders can be found. **Default:** `//medixsrv/Nematodes/data/`

**[-input (path)]**

The path<sup>2</sup> to the directory where the input images can be found. If this argument is specified, "-project" and "-dir" need not be. **Default: Generated from -dir + -project + "input/"**

**[-output (file path)]**

The file path<sup>2</sup> where the output will be written. If this argument is specified, "-project" and "-dir" need not be. **Default: Generated from -dir + -project + "log/feature.log"**

**[-port (number)]**

The port to listen on for node connections. **Default: 8190**

---

<sup>2</sup>Paths must have a trailing "/".

- [-threads (number)]**  
The number of threads to use for network traffic parallelization. **Default: Number of CPU cores**
- [-nodes (number)]**  
The minimum number of nodes that must be joined before beginning the task. The task will automatically begin once this many nodes have connected. **Default: 1**
- [-extension (string)]**  
The extension on the image files. **Default: ".jpg"**
- [-padding (number)]**  
The number of digits in the file name. i.e. "0000001.jpg" has a padding of 7. Legacy video data may have a padding of 6 because I'm bad at conventions. **Default: 7**
- [-frames (number)]**  
The number of frames to process. **Default: File count automatically detected using given parameters. Could be wrong.**

### 3.1.2 Segmentation Parameters

These are all optional, however may need to be configured if segmentation is not successful. Most likely to need tuning is "-threshold\_ratio".

- [-width (number)]**  
The horizontal resolution of the image. **Default: 640**
- [-height (number)]**  
The vertical resolution of the image. **Default: 480**
- [-area\_min (number)]**  
The lower bound for a candidate worm component. If a blob has an area lower than this number, it is discarded and not considered. (e.g. Noise.) **Default: 200**
- [-area\_max (number)]**  
The upper bound for a candidate worm component. If a blob has an area greater than this number, it is discarded and not considered. (e.g. The plate edge.) **Default: 400**
- [-search\_win\_size (number)]**  
The width and height of the crop area. The algorithm crops the frame to this size, centered on the last known worm location, for performance purposes. **Default: 100**
- [-blur\_radius (number)]**  
The width and height of the sliding window used in the box blur. **Default: 3**

**[-threshold\_win\_size (number)]**

The width and height of the sliding window used in the dynamic threshold. This is the neighborhood against which a pixel's intensity will be compared. **Default: 25**

**[-threshold\_ratio (float)]**

The percentage (as a decimal) which a pixel's intensity as a quotient of its neighborhood average must be less than to be considered to be possibly on the worm body. e.g. If this value is 0.9, and the neighborhood average intensity is 255 (pure white), then a pixel's intensity must be below  $0.9 * 255 \approx 229$  in order to be a potential worm body pixel. **Default: 0.9**

### 3.2 Node Arguments

**-mode (host / node)**

This must be specified with "**node**" to start the program in node mode.

**[-host (address)]**

The address of the host process. **Default: "localhost"**

**[-port (number)]**

The port to connect to. **Default: 8190**

**[-threads (number)]**

The number of threads to use for computing. **Default: Number of CPU cores**

### 3.3 JVM Considerations

The JVM heap size is by default 1/4th of your total physical memory. Depending on the input data, a larger heap size may be necessary. You can specify the heap size with the flag **-Xmx#g**, where **#** is some number of gibibytes. For example, your host launch script might read:

```
java -Xmx16g -jar WormSeg.jar -mode host -nodes 8 -project n2_f1
```