Project outline – Most tasks involve creating or updating public artifacts.

| No. | Description | Type | Deadline | Postponements | Status | Est. 5 % | Est. mode | Est. 95 % | Date start | Date end | Actual duration |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **I1** | **Iteration 1** | Iteration | 2019-09-13 | To be determined. | 🟥 | 36:28 | 40:47 | 50:26 | 2019-08-26 | | |
| I1a | Sketch how to adapt SupAmp to RL. → analysis | Task | | 🟨 | ✔ | 3:52 | 5:31 | 6:49 | 2019-08-26 | 2019-08-30 | 13:23 |
| I1b | Sketch how to model supervisor failures. | Task | | | 🟥 | | | | | | 19h 40m |
| | - Add my insight that I'm not just writing for myself to supamp-reamp. And that the analysis is preliminary. | | | | | | | | | | |
| I1c | Create an empty Draft Basis and fill in as far as possible. | Task | | | | | | | | | |
| I1d | Announce my project on LW or MxD. | Task | | | | | | | | | |
| | Announce search for writing partner on LW or MxD. | Task | | | | | | | | | |
| | Paul's code for SupAmp runs on my machine and I roughly know my way around it. | Task | | | | | | | | | |
| | Read and summarize relevant literature. | Task | | | | | | | | | |
| | **Iteration 2** | Iteration | | | | | | | | | |
| | Study missing ML basics. – ML, deep learning, RL | Task | | | | | | | | | |
| | Verify design so far. | Task | | | | | | | | | |
| | Design how to adapt SupAmp to RL. | Task | | | | | | | | | |
| | Fill in Draft Basis further. | Task | | | | | | | | | |
| | Hopefully found writing partner(s). | Task | | | | | | | | | |
| | **Iteration 3** | Iteration | | | | | | | | | |
| | Adapted SupAmp to RL. | Task | | | | | | | | | |
| | Run some experiments from CSASupAmp with RL instead of SL. | Task | | | | | | | | | |
| | Write short article about the differences between SupAmp and ReAmp. | Task | | | | | | | | | |
| | **Iteration 4** | Iteration | | | | | | | | | |
| | Design experiments for ReAmp with overseer failures. | Task | | | | | | | | | |
| | Design changes to ReAmp to accommodate experiments. | Task | | | | | | | | | |
| | **Iteration 5** | Iteration | | | | | | | | | |
| | Adapt ReAmp code. | Task | | | | | | | | | |
| | Run experiments. | Task | | | | | | | | | |
| | Finish filling in Draft Basis. | Task | | | | | | | | | |
| | **Iteration 6** | Iteration | | | | | | | | | |
| | Revisit literature. | Task | | | | | | | | | |
| | Make writing plan. | Task | | | | | | | | | |
| | Make build pipeline for article. | Task | | | | | | | | | |
| | **Iteration 7** | Iteration | | | | | | | | | |
| | Write draft. | Task | | | | | | | | | |
| | Revise draft. | Task | | | | | | | | | |
| | Solicit feedback. | Task | | | | | | | | | |
| | **Iteration 8** | Iteration | | | | | | | | | |
| | Write final version. | Task | | | | | | | | | |
| | Submit article. | Task | | | | | | | | | |

## Abbreviations/Glossary/Bibliography

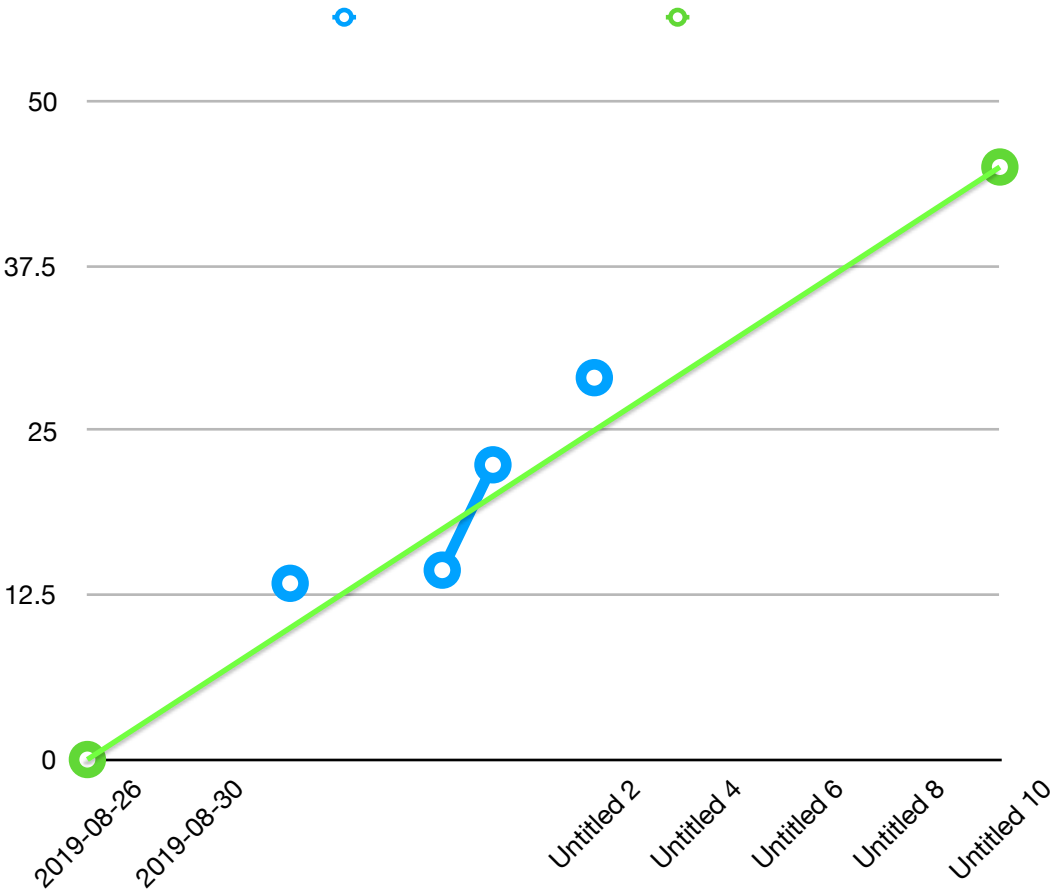| CoR | Booth et al.: The Craft of Research |
|---|---|
| CSASupAmp | Christiano et al.: Supervising strong learners by amplifying weak experts |
| Est. 5 % | 5th percentile of my estimated duration distribution/leftmost point in triangle distribution |
| Est. mode | mode of my estimated duration distribution |
| Est. 95 % | 95th percentile of my estimated duration distribution/rightmost point in triangle distribution |
| Draft Basis | A template derived from CoR, p. 175, which when filled in completely, provides all the information necessary for planning a draft. Includes the structure of the argument. |
| LW | LessWrong |
| MxD | MIRIxDiscord |
| RL | reinforcement learning |
| ReAmp | SupAmp adapted to RL |
| SL | supervised learning |
| SupAmp | The system for iterated distillation and amplification using supervised learning from CSASupAmp |
| | |
| | |
| | |
| | |

## Table 1

| | | | |
|---|---|---|---|
| 2019-08-26 | 0 | | |
| 2019-08-27 | | | |
| 2019-08-28 | | | |
| 2019-08-29 | | | |
| 2019-08-30 | | 13.3833333333333 | 13:23 |
| 2019-08-31 | | | |
| 2019-09-01 | | | |
| 2019-09-02 | | 14.3833333333333 | |
| 2019-09-03 | | 22.3833333333333 | |
| 2019-09-04 | | | |
| 2019-09-05 | | 29 | |
| 2019-09-06 | | | |
| 2019-09-07 | | | |
| 2019-09-08 | | | |
| 2019-09-09 | | | |
| 2019-09-10 | | | |
| 2019-09-11 | | | |
| 2019-09-12 | | | |
| 2019-09-13 | 45 | | |

## Table 2

| | | | |
|---|---|---|---|
| 2019-08-30 | | 13.3833333333333 | 13:23 |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |

This thing is wrong.

## Estimates preprocessed for Dugless

| | Est. mode | Est. 5 % | Est. 95 % | Est. mode ratio |
|---|---|---|---|---|
| | 2447 | 2188 | 3026 | 0.309069212410501 |
| | 331 | 232 | 409 | 0.559322033898305 |
| | 120 | 60 | 180 | 0.5 |
| | 180 | 90 | 360 | 0.333333333333333 |
| | 60 | 30 | 90 | 0.5 |
| | 15 | 10 | 30 | 0.25 |
| | 120 | 30 | 360 | 0.272727272727273 |
| | 60 | 15 | 360 | 0.130434782608696 |
| | 30 | 15 | 90 | 0.2 |
| | 30 | 15 | 60 | 0.333333333333333 |
| | 30 | 15 | 60 | 0.333333333333333 |
| | 30 | 15 | 60 | 0.333333333333333 |
| | 120 | 60 | 360 | 0.2 |
| | 30 | 15 | 60 | 0.333333333333333 |
| | 120 | 0 | 360 | 0.333333333333333 |
| | 30 | 15 | 60 | 0.333333333333333 |
| | 90 | 30 | 240 | 0.285714285714286 |
| | 240 | 90 | 480 | 0.384615384615385 |
| | 780 | 540 | 1500 | 0.25 |
| | 90 | 45 | 150 | 0.428571428571429 |