

Optimizing labeling

The limits of weakly supervised learning in osteophyte severity grading and localisation in 4 hip quadrants



Severity grade 0 through 3 of superior femoral osteophyte based on the OARSI severity scale

AUTHORS

Alvin Ye¹, a.d.ye@student.tudelft.nl
Supervisors: Gijs van Tulder¹, Jesse Krijthe¹

AFFILIATIONS

Delft University of Technology¹

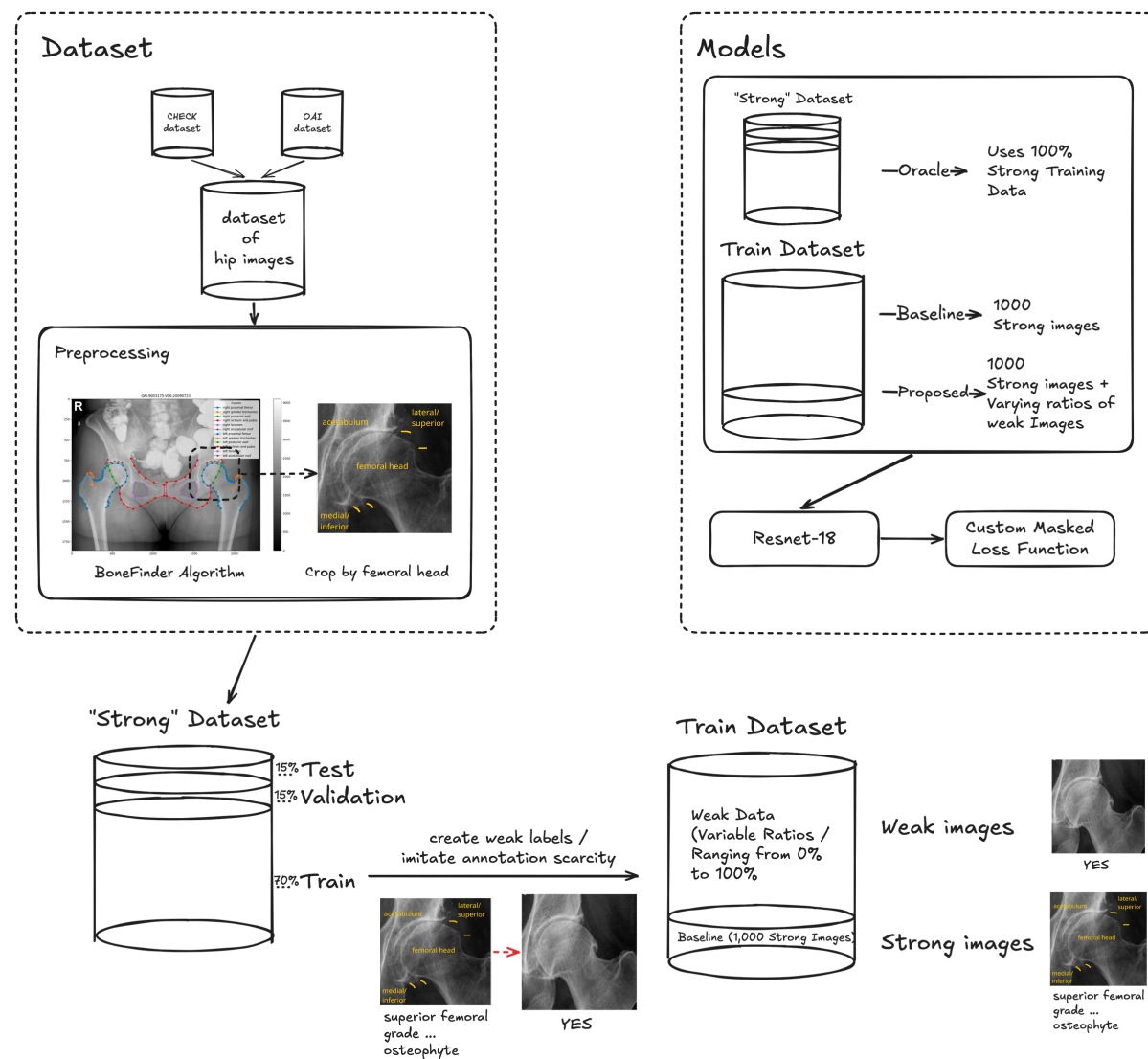
1. THE 500 MILLION PROBLEM

- Osteoarthritis (OA) is a common disease that develops osteophytes.
- In 2020, about 595 million people had a form of osteoarthritis.
- Detailed osteophyte annotation in X-Ray Images is time-consuming and costly requiring an expert.
- Consequently, clinical datasets are expensive to scale with annotated data.
- We categorize weak data with a binary indicator of osteophytes in the entire image. We categorize strong data as expert annotated severity grade in the 4 hip joint locations
- It is unknown to what extent adding weak data continues to help grading severity and localising osteophytes.

THE MAIN QUESTION

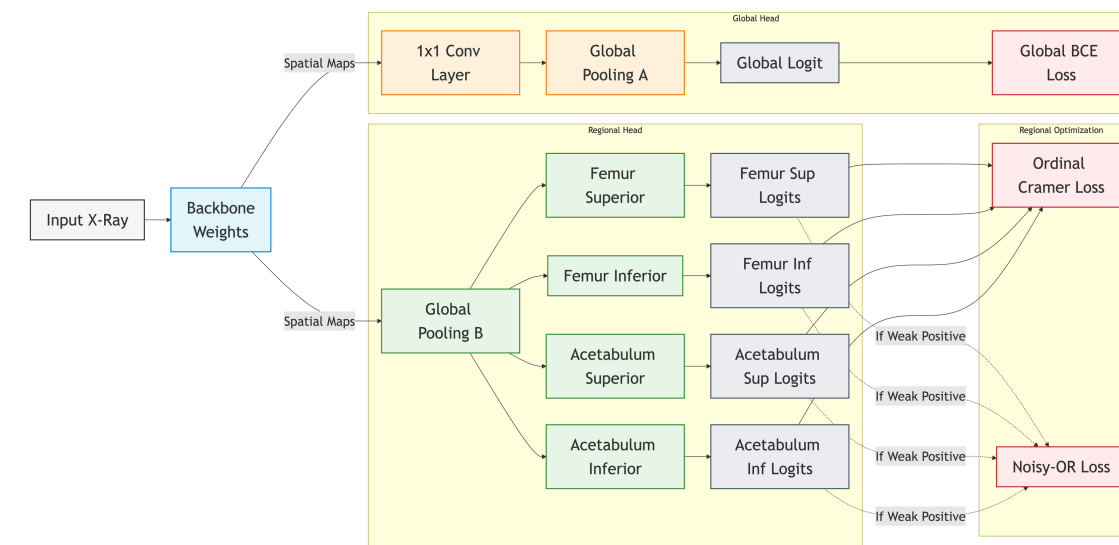
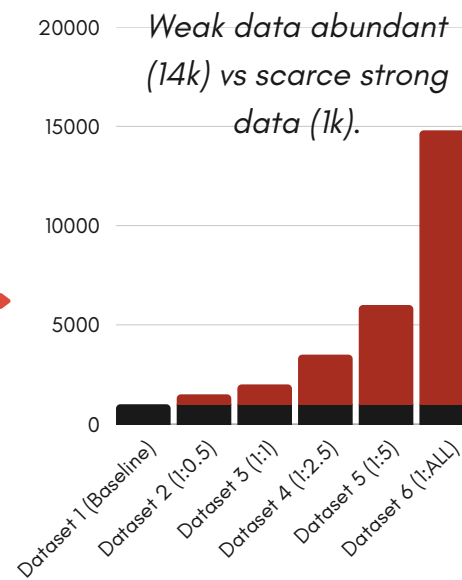
“ How do different weak-to-strong data mixes change performance in a multi task machine learning pipeline and specifically when does it stop improving?

2. TRAINING DATA



Preprocessing the datasets to create a Train/Validation/Test purely on strong images. The Train split is then modified to become 1000 strong images + remaining becomes weak.

3. FOUR CONFIGURATIONS



Two headed architecture to handle two types of data

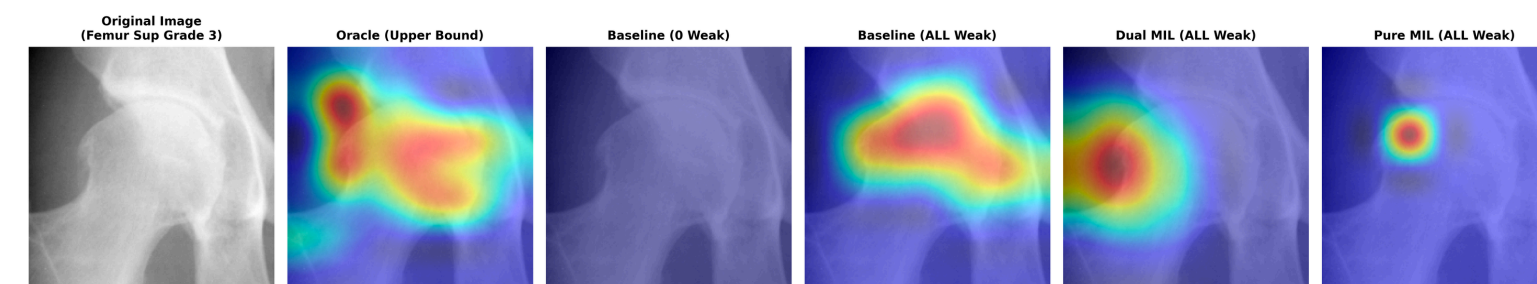
We evaluate 4 model configurations across all six datasets:

- Oracle Model (Upper Bound): Trained on 100% strong data using the dual-head pipeline.
- Masked Baseline: Dual-head pipeline, but ignores weak positives (uses [0,0,0,0] ground truth for healthy joints).
- Dual-Head MIL (Proposed): Full architecture capabilities, leveraging Noisy-OR Loss.
- Pure MIL Ablation: Regional head only (isolates and verifies the dual-head benefit).
- Quadratic Weighted Kappa as metric indicating agreement with experts (0 = random 1 = perfect agreement). QWK score is lower when predictions fall further from true grade.

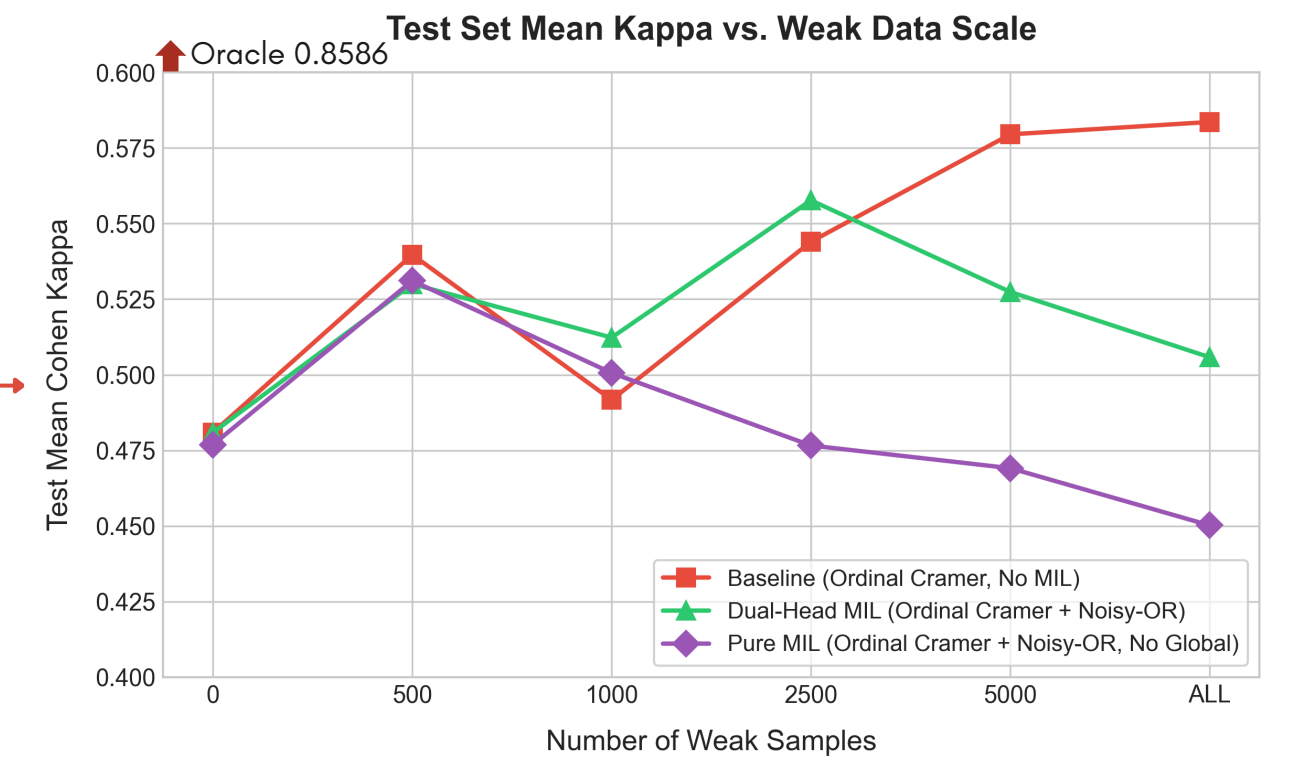
4. RESULTS & KEY FINDINGS

Model	ρ (Weak Samples)	Test Mean Kappa
Oracle	100% Strong	0.8586
Masked Baseline	0	0.4809
Masked Baseline	ALL (~15k) (Peak)	0.5836
Dual-Head MIL	0	0.4809
Dual-Head MIL	1000 (Val-Selected Peak)	0.5123
Dual-Head MIL	2500 (Test-Absolute Peak)	0.5577
Dual-Head MIL	ALL (~15k)	0.5058
Pure MIL	0	0.4769
Pure MIL	500 (Peak)	0.5311
Pure MIL	ALL (~15k)	0.4502

Masked baseline scales as data scales to the maximum (Kappa 0.48 \rightarrow 0.58). MIL peak early then degrade. Validation selected model does not contain the highest test kappa.



Grading-vs-localisation trade-off: masked baseline grades best but shows diffusely like the Oracle, Dual-MIL localises tightly but grades worse. Pure MIL misses the correct area



Scaling data ratio shows that Pure MIL Kappa peaks early and consistently drops, whereas Dual-Head peaks later but also subsequently drops. Masked Baseline improves but sees an anomaly at 1000 with the other networks

	Baseline (0 Weak) - Lower Bound	Baseline (ALL Weak) - Best Baseline	Dual MIL (ALL Weak) - Best MIL
Grade 0	2154 (89.4%)	2166 (89.9%)	1930 (80.1%)
Grade 1	209 (46.3%)	189 (37.5%)	111 (24.6%)
Grade 2	96 (34.9%)	55 (20.0%)	48 (17.5%)
Grade 3	8 (30.8%)	5 (19.2%)	5 (19.2%)
Grade 0	217 (9.0%)	194 (8.1%)	423 (17.6%)
Grade 1	51 (11.3%)	211 (46.8%)	269 (59.6%)
Grade 2	81 (29.5%)	117 (42.5%)	114 (41.5%)
Grade 3	3 (1.1%)	99 (36.0%)	109 (39.6%)
Grade 0	1 (0.0%)	1 (0.0%)	52 (2.2%)
Grade 1	0 (0.0%)	2 (0.4%)	66 (14.6%)
Grade 2	0 (0.0%)	4 (1.5%)	66 (14.6%)
Grade 3	3 (11.5%)	3 (11.5%)	4 (1.5%)
Grade 0	1864 (77.4%)	2329 (96.7%)	2409 (100.0%)
Grade 1	127 (28.2%)	56 (12.4%)	0 (0.0%)
Grade 2	44 (16.0%)	25 (9.3%)	0 (0.0%)
Grade 3	5 (19.2%)	1 (3.8%)	0 (0.0%)
Grade 0	538 (22.3%)	59 (2.4%)	0 (0.0%)
Grade 1	311 (69.0%)	372 (82.5%)	451 (100.0%)
Grade 2	196 (71.3%)	30 (10.4%)	0 (0.0%)
Grade 3	5 (19.2%)	7 (26.9%)	0 (0.0%)
Grade 0	7 (0.3%)	20 (0.8%)	0 (0.0%)
Grade 1	12 (2.7%)	23 (5.1%)	0 (0.0%)
Grade 2	35 (12.7%)	213 (77.5%)	275 (100.0%)
Grade 3	0 (0.0%)	18 (69.2%)	0 (0.0%)
Grade 0	0 (0.0%)	0 (0.0%)	0 (0.0%)
Grade 1	0 (0.0%)	0 (0.0%)	0 (0.0%)
Grade 2	0 (0.0%)	0 (0.0%)	0 (0.0%)
Grade 3	0 (0.0%)	26 (100.0%)	0 (0.0%)

The gain in Kappa score comes from clean negatives sharpening the Grade 0/1 boundary not from learning osteophytes (masked never sees a weak positive). MIL models struggle with predicting the correct grade due to missing severity grading information beyond Grade 1

5. LIMITATIONS

- Dataset is skewed, only few severe cases to train on
- Single seed
- Best model ignores weak positive leaving good data unused
- Qualitative GradCAM instead of quantitative

6. CONCLUSION & FUTURE WORK

- Weak labels provide better prediction, only when using healthy samples
- Best grading performance is not always best localisation performance
- Missing severity labels provide noisy data dropping performance
- Vision-language extraction of rich features (severity) may be a solution
- Future work can constrain spatial analysis using anatomical priors
- Future work should verify results statistically