# aival

Healthcare AI evaluation and monitoring

# AI will transform healthcare

Aival enables an AI product to be properly analysed, to give clinicians the confidence to know that it will work at their local site and for their patients

***Our solution allows rapid, scalable and repeatable independent assessment of AI products without requiring technical expertise***

Our methodology is based on decades of experience in developing medical imaging AI algorithms and commercial products, understanding their failure modes and weaknesses and how to test for them

# Our founder

has 15+ years' experience in research, development, and regulation of AI products for healthcare

## Kanwal Bhatia, Ph.D.

Head of Data Science at Visulytix, leading a team of 6 data scientists. Developed IP that sold to big pharma / device manufacturers

AI Architect at Odin Vision

Technical Advisory Board at Ultromics

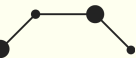Ph.D. in Medical Image Computing (Imperial College London, 2007) with 1500+ citations
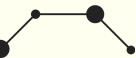
# Clinical adoption of AI is slow

- **> 500 AI devices cleared for clinical use (FDA)**
  - $5bn investment into medical imaging AI since 2015

- **Lack of standard pathways to adoption**
  - AI products are hard to understand, operating as 'black boxes'
  - No standard pathways for validating products before adoption (current methods are expensive / biased)
  - Weak monitoring of AI performance once in use
  - Clinical staff do not have the time or skillset to evaluate technical performance and safety of AI
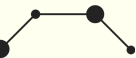
*\* Signify Research, May 2023*

# Build trust through evaluation

Aival software evaluates AI products rapidly and at scale

- ## Comparison
  - Identify and compare products that provide greatest clinical benefit for a given site

- ## Evaluation
  - Rigorous *independent* assessment
  - Substantiate manufacturer performance claims
  - Accelerate time to sale and reduce cost of adoption

- ## Monitoring
  - Ensure products continue to perform as expected over time
  - Standardise post-market surveillance reporting

# AI product assessment on local site data

### Performance metrics

How well does it work?
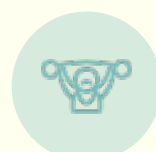Does it work across all acquisition devices and pathologies

### Explainability

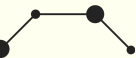Did the algorithm make the right prediction for the right reasons?

### Fairness & bias

Are all population subgroups treated in the same way?

### Robustness

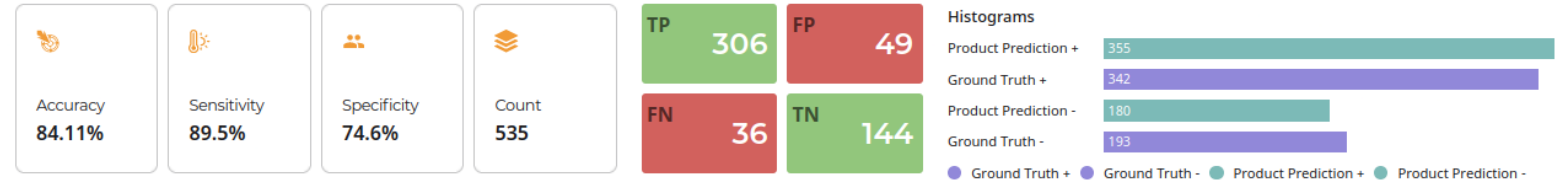Will the algorithm perform just as well with unexpected / variable data?
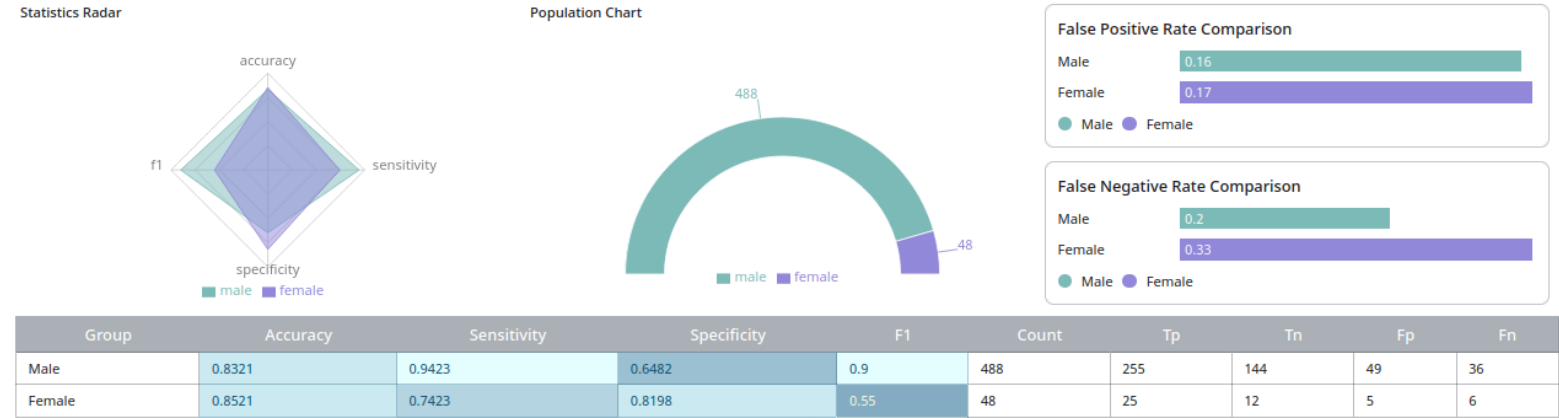
# Performance

Accuracy
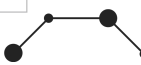
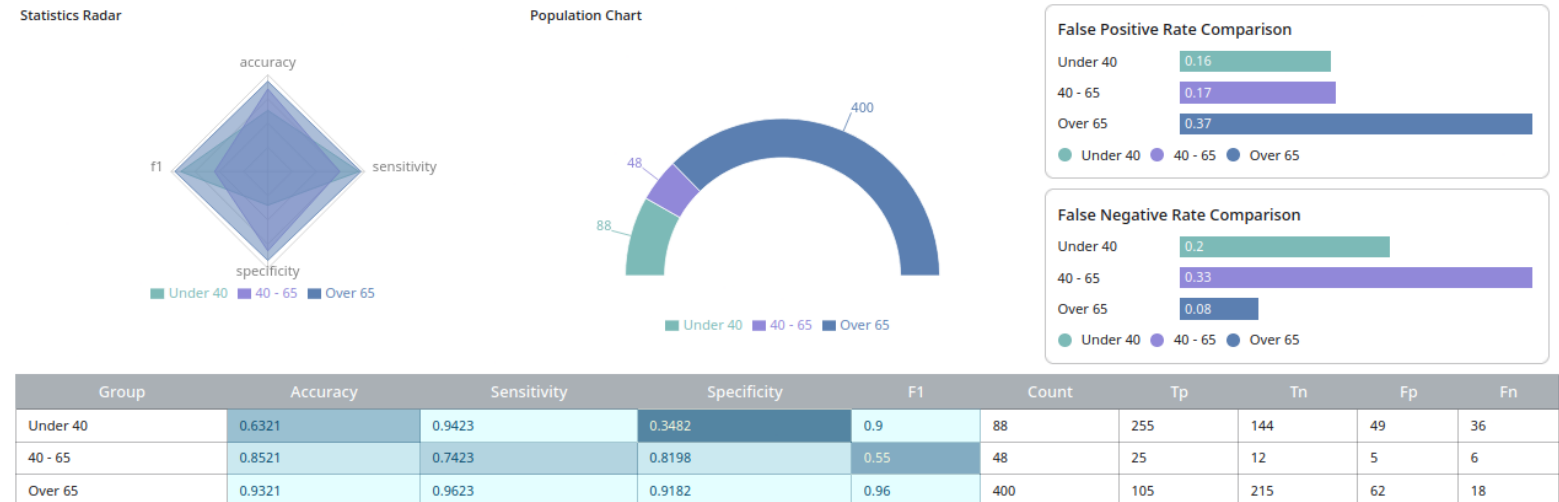Sensitivity

Specificity

Errors

Subgroup analysis

## Overall Evaluation Results

| Accuracy | Sensitivity | Specificity | Count |
|----------|-------------|-------------|-------|
| 84.11% | 89.5% | 74.6% | 535 |

| TP | 306 | FP | 49 |
|----|-----|----|----|
| FN | 36 | TN | 144 |

### Histograms

| | |
|---|---|
| Product Prediction + | 355 |
| Ground Truth + | 342 |
| Product Prediction - | 180 |
| Ground Truth - | 193 |

● Ground Truth +  ● Ground Truth -  ● Product Prediction +  ● Product Prediction -

## Sex Results Breakdown

Statistics Radar



● male  ● female

Population Chart

488

48

● male  ● female

### False Positive Rate Comparison

| Male | 0.16 |
|------|------|
| Female | 0.17 |

● Male  ● Female

### False Negative Rate Comparison

| Male | 0.2 |
|------|-----|
| Female | 0.33 |

● Male  ● Female

| Group | Accuracy | Sensitivity | Specificity | F1 | Count | Tp | Tn | Fp | Fn |
|-------|----------|-------------|-------------|-----|-------|-----|-----|-----|-----|
| Male | 0.8321 | 0.9423 | 0.6482 | 0.9 | 488 | 255 | 144 | 49 | 36 |
| Female | 0.8521 | 0.7423 | 0.8198 | 0.55 | 48 | 25 | 12 | 5 | 6 |

## Age Results Breakdown

Statistics Radar



● Under 40  ● 40 - 65  ● Over 65

Population Chart

400

48

88

● Under 40  ● 40 - 65  ● Over 65

### False Positive Rate Comparison

| Under 40 | 0.16 |
|----------|------|
| 40 - 65 | 0.17 |
| Over 65 | 0.37 |

● Under 40  ● 40 - 65  ● Over 65

### False Negative Rate Comparison

| Under 40 | 0.2 |
|----------|-----|
| 40 - 65 | 0.33 |
| Over 65 | 0.08 |

● Under 40  ● 40 - 65  ● Over 65

| Group | Accuracy | Sensitivity | Specificity | F1 | Count | Tp | Tn | Fp | Fn |
|-------|----------|-------------|-------------|-----|-------|-----|-----|-----|-----|
| Under 40 | 0.6321 | 0.9423 | 0.3482 | 0.9 | 88 | 255 | 144 | 49 | 36 |
| 40 - 65 | 0.8521 | 0.7423 | 0.8198 | 0.55 | 48 | 25 | 12 | 5 | 6 |
| Over 65 | 0.9321 | 0.9623 | 0.9182 | 0.96 | 400 | 105 | 215 | 62 | 18 |

# Fairness

Accuracy parity

Predictive value parities

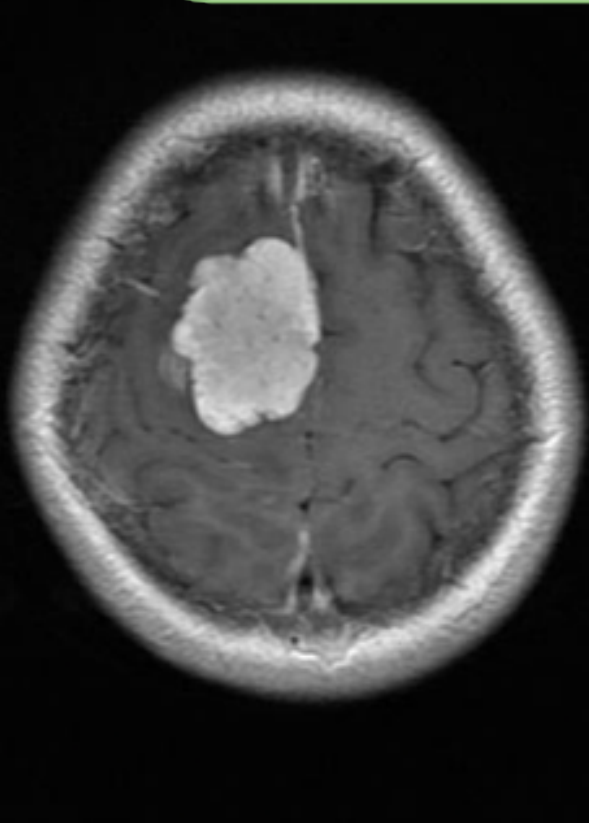False positive rate parity
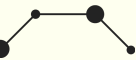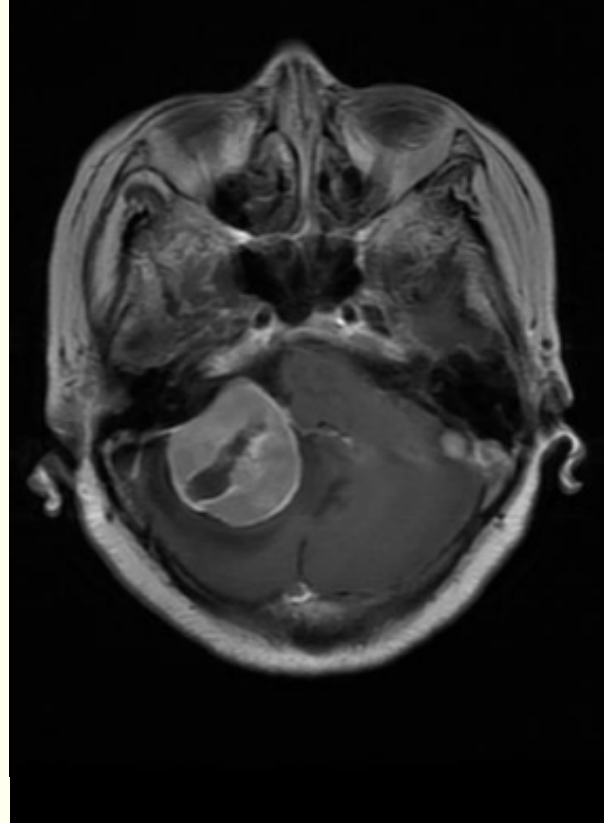
False negative rate parity

# Explainability (black-box)

Is the AI product making the right decisions for the right reasons? We test products as black-boxes without access to underlying model / architecture



Classification: Meningioma, confidence=0.999

Classification: Meningioma, confidence=0.998

# Explainability (black-box)

Is the AI product making the right decisions for the right reasons? We test products as black-boxes without access to underlying model / architecture
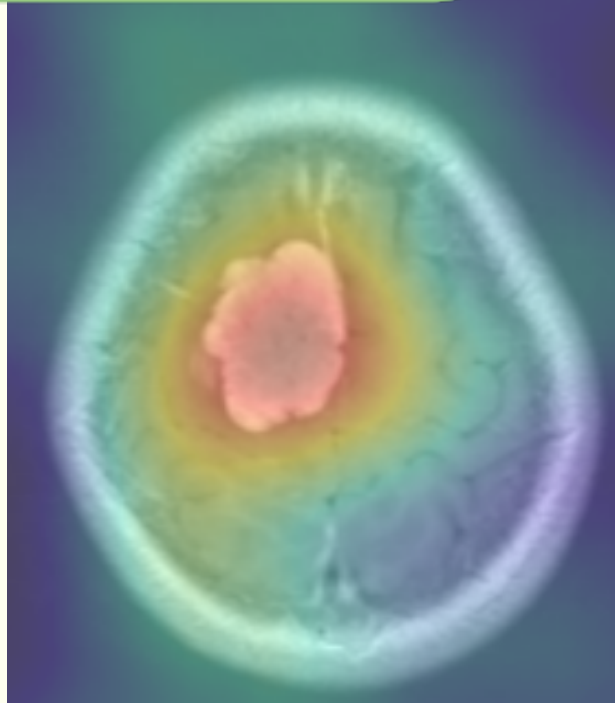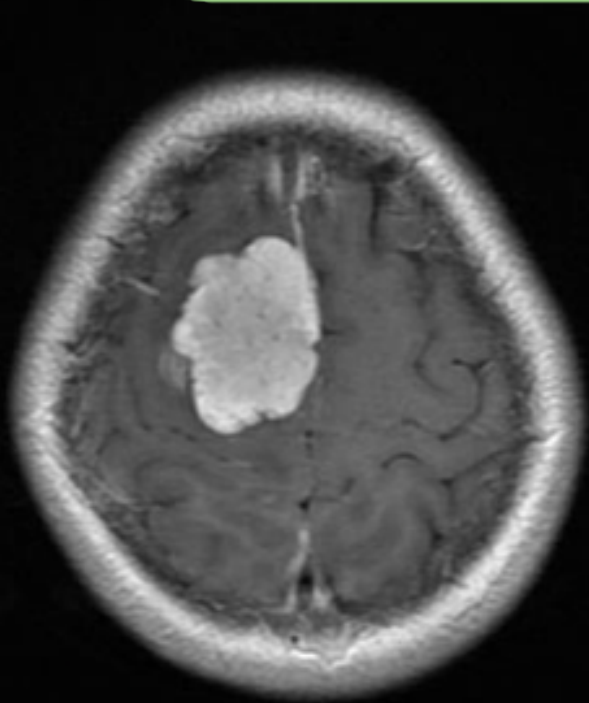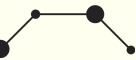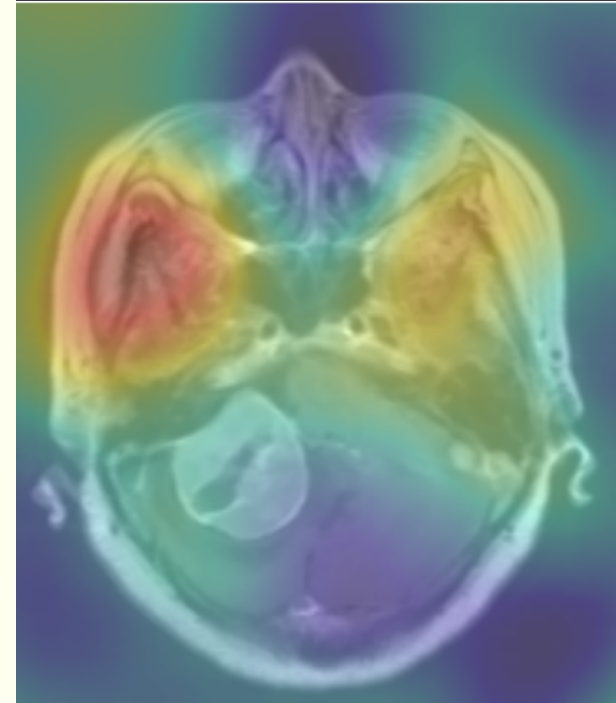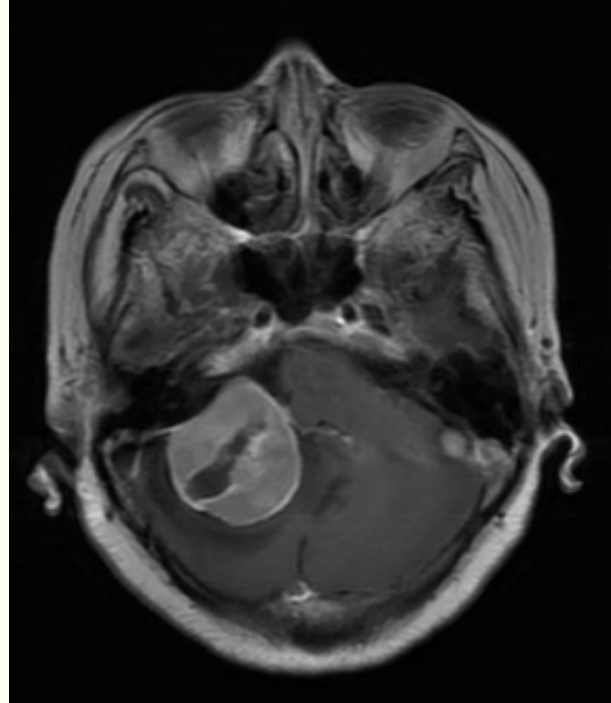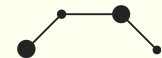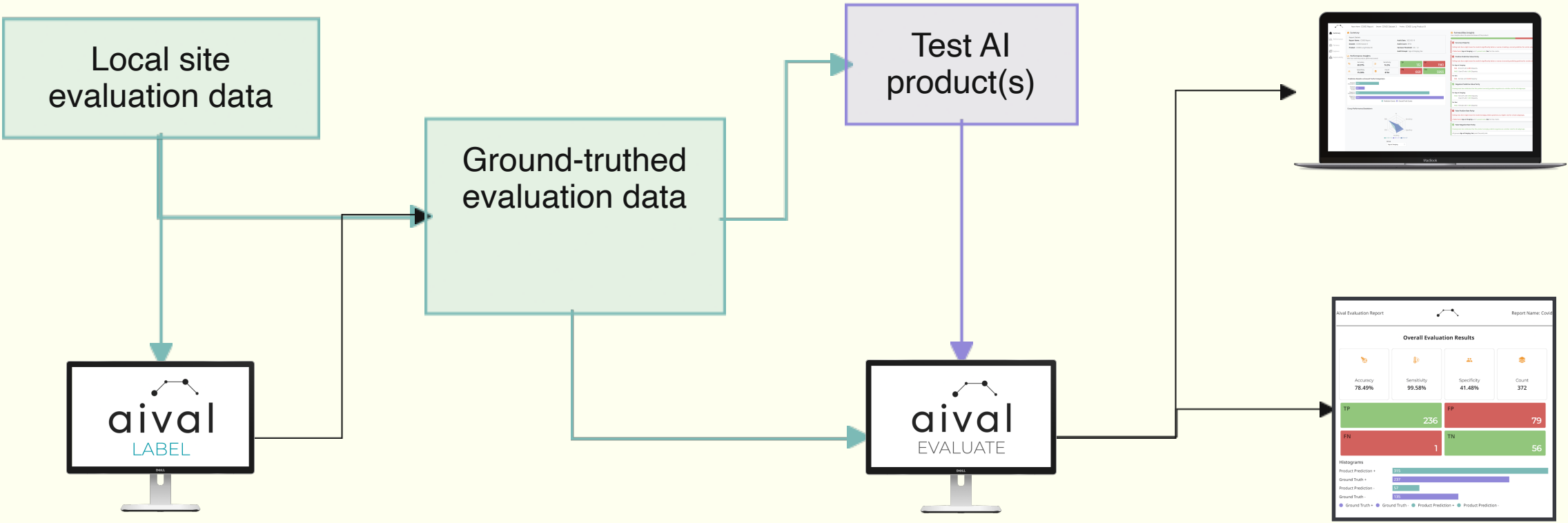


Classification: Meningioma, confidence=0.999

Classification: Meningioma, confidence=0.998

# Aival evaluation workflow

# Sample analysis report



**Report Name** COVID Report    **Dataset** COVID Dataset 3    **Product** COVID Lung Product B

- Summary
- Performance
- Fairness
- Explorer
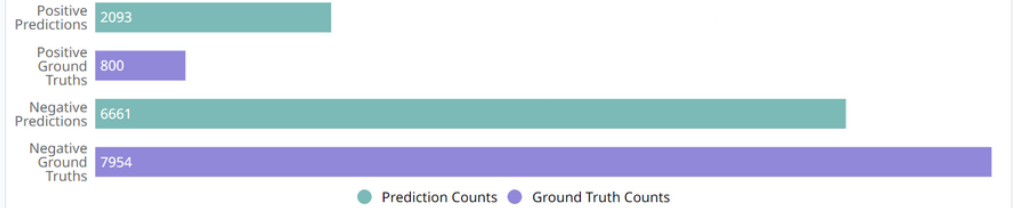- Explainability

## 🏠 Summary

### Report Details

| | |
|---|---|
| **Report Name:** COVID Report | **Audit Date:** 2023-05-19 |
| **Dataset:** COVID Dataset 3 | **Audit Count:** 8754 |
| **Product:** COVID Lung Product B | **Fairness Threshold:** 0.8 - 1.2 |
| | **Audit Groups:** Age at Imaging, Sex |

## 📊 Performance Insights
View how well this product performed overall.

| Accuracy 69.97% | Sensitivity 16.5% | TP 132 | FP 1961 |
|---|---|---|---|
| Specificity 75.35% | Count 8754 | FN 668 | TN 5993 |

### Prediction Results vs Ground Truths Comparison

- Positive Predictions 2093
- Positive Ground Truths 800
- Negative Predictions 6661
- Negative Ground Truths 7954

● Prediction Counts  ● Ground Truth Counts

### Group Performance Breakdown

F1 / FNR / Sensitivity / FPR / Specificity

## ⚖️ Fairness/Bias Insights
View insights about the potential biases of this product.

9 / 15 passed tests

### ❌ Accuracy Disparity

Failing tests here might mean the model is significantly better or worse at making a correct prediction for certain subgroups.

2 failed tests (**Age at Imaging**) and 1 passed tests (**Sex**) for this metric.

### ❌ Positive Predictive Value Parity

Failing tests here might mean the model is significantly better or worse at correctly predicting positives for certain subgroups than in others.

**For Age at Imaging:**                        Reference Group: **Under 50**
FAIL  50 to 67 with **2.09X** Disparity
PASS  Over 67 with **1.03X** Disparity

**For Sex:**                                  Reference Group: **Male**
FAIL  Female with **0.62X** Disparity

### ✅ Negative Predictive Value Parity

Passing tests here indicates that the product correctly predicts negatives at a similar rate for all subgroups.

**For Age at Imaging:**                        Reference Group: **Under 50**
PASS  50 to 67 with **0.98X** Disparity
      Over 67 with **1.02X** Disparity

**For Sex:**                                  Reference Group: **Male**
PASS  Female with **1.03X** Disparity

### ❌ False Positive Rate Parity

Failing tests here might mean the model wrongly predicts positives at a higher rate for certain subgroups.

2 failed tests (**Age at Imaging**) and 1 passed tests (**Sex**) for this metric.

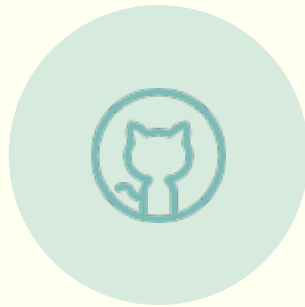### ✅ False Negative Rate Parity

Passing tests here indicates that the product wrongly predicts negatives at a similar rate for all subgroups.

# Use cases

### Healthcare Providers

- Validate manufacturer performance claims on local data
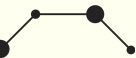- Ensure fairness across demographics
- Compare different AI products

### AI Vendors

- Gain trust with clinical users
- Internal self-assessment of failure modes
- Standardise reporting for regulatory submissions

### AI Platforms

- Help your users to assess different products across your platform

# Get in touch

@ kanwal@aival.io

☎ +447795975256

🌐 https://www.aival.io

aival