

Early identification of polycystic ovary syndrome and associated multimorbidity — a multi-cohort study

Nouman Ahmed¹, Malgorzata Wamil^{1,2,3}, Nathalie Conrad^{1,4}, Arash Mohazzab⁵, Jie Lian¹, Kazem Rahimi^{1*}, Shishir Rao^{1†*}

¹ Nuffield Dept. of Women's & Reproductive Health, University of Oxford, UK
² Mayo Clinic Healthcare, London, UK
³ Great Western Hospital NHS Trust, Swindon, UK

⁴ Dept. of Cardiovascular Sciences, KU Leuven, Belgium
⁵ School of Public Health, Iran University of Medical Sciences, Tehran, Iran
[†] Equal contribution · ^{*} Corresponding authors

01 BACKGROUND

A common syndrome, frequently missed.

PCOS affects 5–18% of women of reproductive age, yet up to 70% never receive a formal diagnosis.

PCOS is a heterogeneous, multisystem condition emerging gradually across adolescence and early adulthood — defined by oligo-/anovulation, hyperandrogenism, and polycystic ovarian morphology (Rotterdam, 2003).

Why earlier matters

- 4× increased risk of type 2 diabetes
- Elevated CV risk — hypertension, dyslipidaemia, IHD, stroke
- MASLD, endometrial hyperplasia & cancer
- Adverse obstetric outcomes; high rates of anxiety & depression
- Conventional CV risk scores systematically underestimate risk in younger women

02 AIM & APPROACH

Predict PCOS years before recognition.

Develop and externally validate **TRisk**, a transformer-based survival model on routinely collected EHR, to predict future clinical PCOS diagnosis and identify women at elevated risk of PCOS-related multimorbidity.

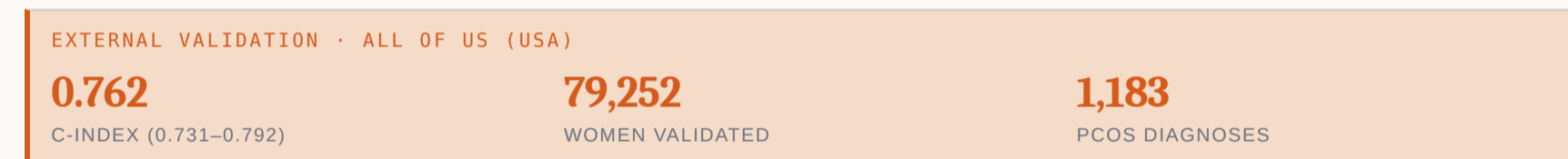
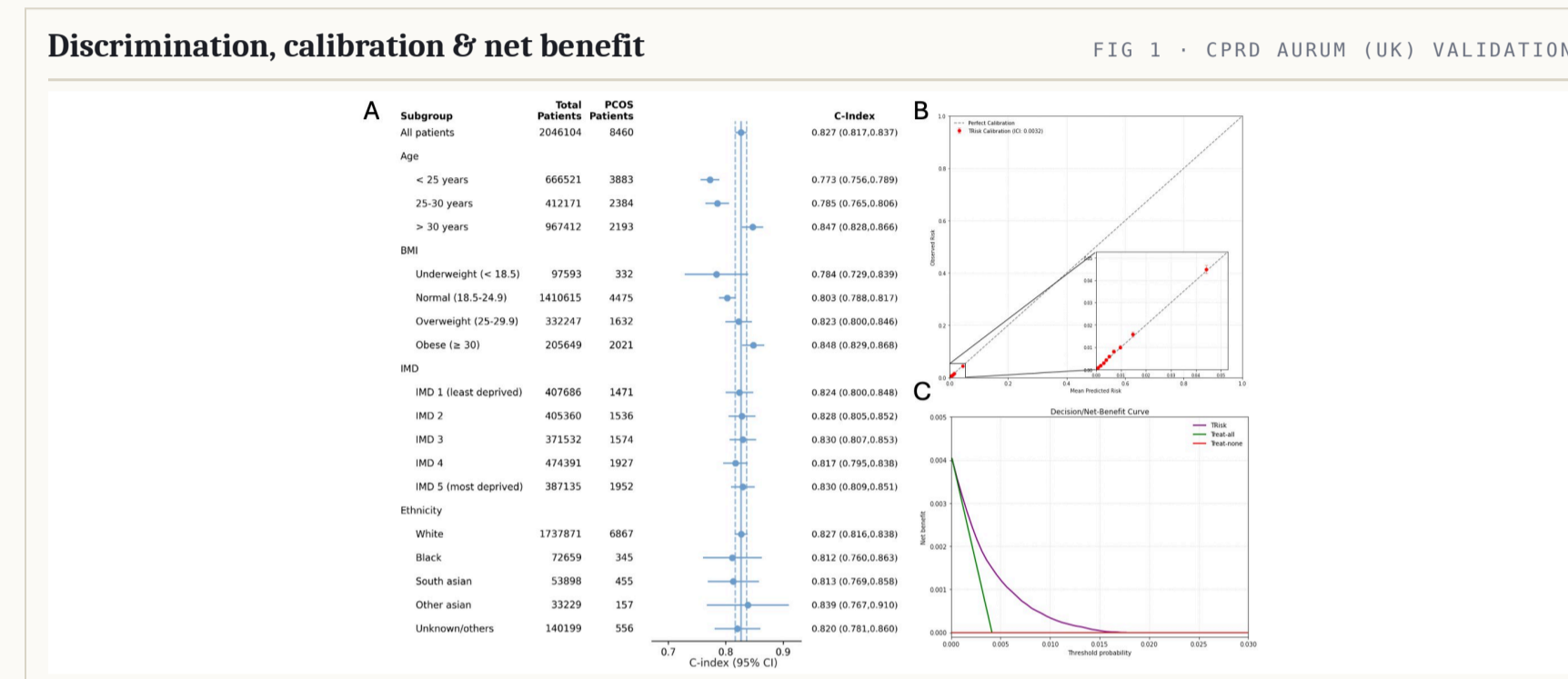
03 COHORTS

UK development. USA external validation.

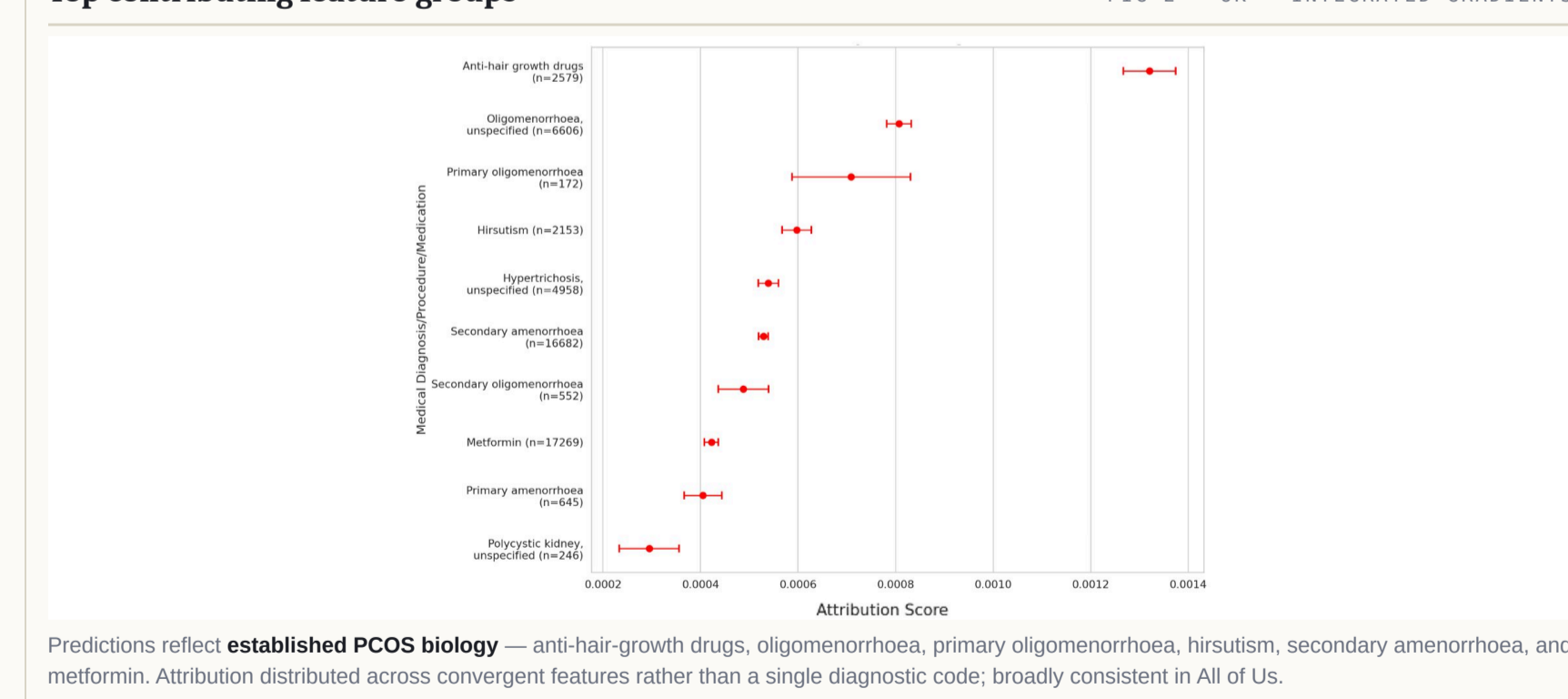
UK · INTERNAL		USA · EXTERNAL	
CPRD Aurum		All of Us	
DEVELOPMENT & VALIDATION			
6,507,815	2,046,104	52,834	79,252
DEVELOPMENT		FINE-TUNE	
37,805	1.5 yr	1,966	2.67 yr
TOTAL PCOS PATIENTS		TOTAL PCOS PATIENTS	
	MEDIAN F/U		MEDIAN F/U

04 RESULTS · MODEL PERFORMANCE

Strong, stable discrimination across systems.



Top contributing feature groups



INTERPRETATION

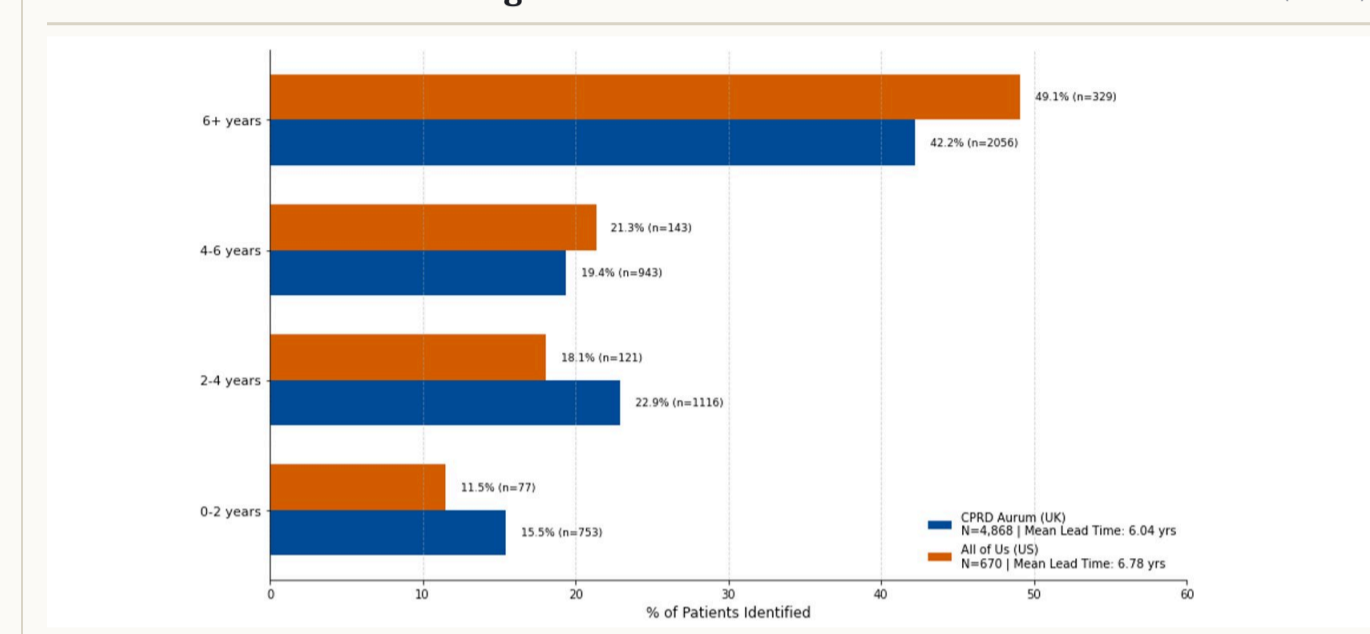
A transformer-based survival model on routinely collected EHR can identify women at elevated PCOS risk several years before clinical recognition — and stratify a high-risk group with substantial, modifiable multimorbidity burden currently missed by existing care pathways. Findings are consistent across the UK (15-yr follow-up) and USA (10-yr follow-up) healthcare systems.

IMPLICATIONS

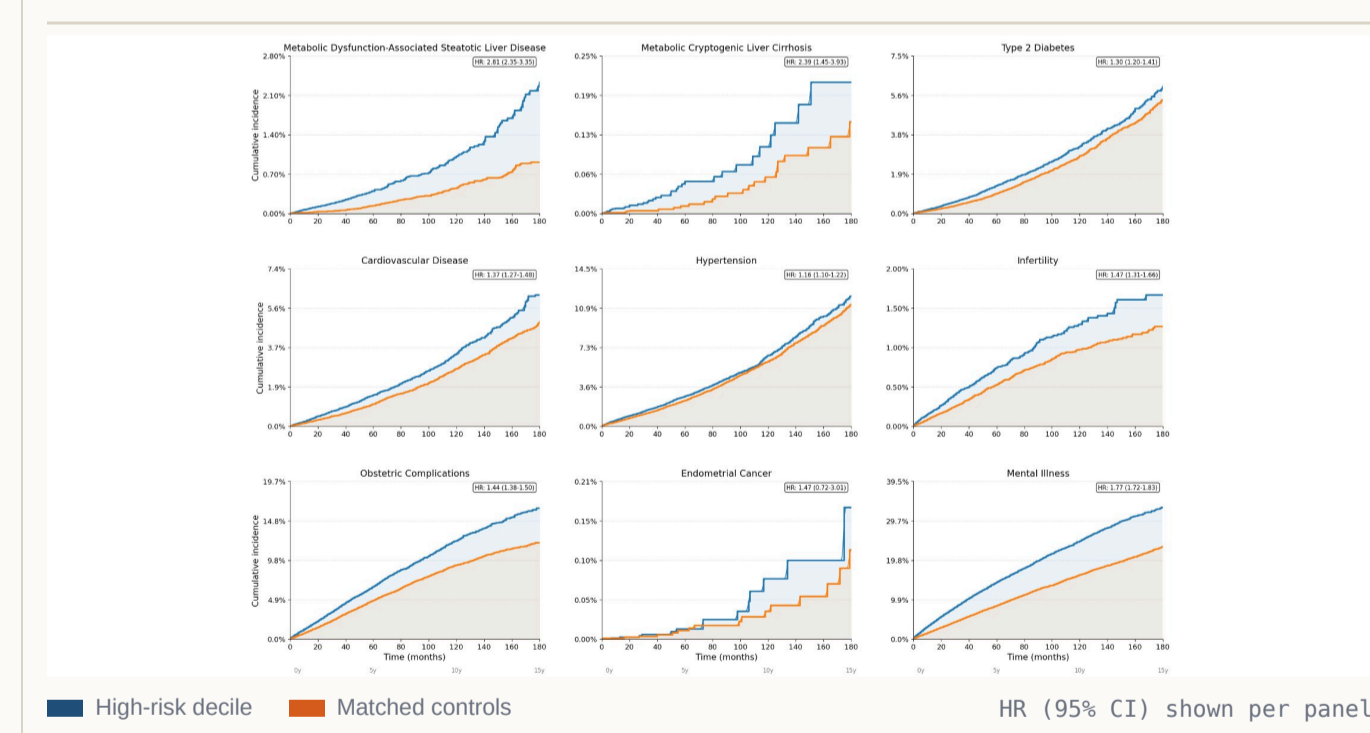
Deployed as a background risk-flagging tool in primary care, this approach could support structured case-finding, earlier metabolic and CV monitoring, and timely fertility planning. Prospective evaluation is needed to confirm earlier identification reduces diagnostic delay and improves outcomes — alongside equity, clinical acceptability, and cost-effectiveness.

05 LEAD TIME & HIGH-RISK PHENOTYPE ANALYSIS

Lead time before clinical diagnosis



Cumulative incidence — top decile vs. matched controls



ALL OF US (USA) · 10-YEAR MATCHED COHORT FOLLOW-UP

Broadly consistent with UK findings across cardiometabolic, reproductive, and mental health outcomes.

Cardiometabolic	Reproductive	Mental health
MASLD · T2D · CVD · HTN	Infertility · obstetric	Consistent with UK finding

LIMITATIONS

EHR-based outcomes capture recorded diagnoses — undiagnosed cases remain censored. UK-trained model required fine-tuning before USA validation. Discrimination lower in younger and underweight subgroups. Calibration deviated mildly at high predicted probabilities in the USA cohort.

CONTACT

nouman.ahmed@wrh.ox.ac.uk · Deep Medicine · University of Oxford