

Time-independent disease state identification defines distinct trajectories determined by localised vs systemic inflammation in patients with early rheumatoid arthritis

Steinz, Nils; Maarseveen, Tjardo D.; van den Akker, Erik B.; Cope, Andrew P.; Isaacs, John D.; Winkler, Aaron R.; Huizinga, Tom W.J; Abraham, Yann; Knevel, Rachel

DOI

[10.1016/j.ard.2025.04.011](https://doi.org/10.1016/j.ard.2025.04.011)

Publication date

2025

Document Version

Final published version

Published in

Annals of the Rheumatic Diseases

Citation (APA)

Steinz, N., Maarseveen, T. D., van den Akker, E. B., Cope, A. P., Isaacs, J. D., Winkler, A. R., Huizinga, T. W. J., Abraham, Y., & Knevel, R. (2025). Time-independent disease state identification defines distinct trajectories determined by localised vs systemic inflammation in patients with early rheumatoid arthritis. *Annals of the Rheumatic Diseases*, 84(8), 1301-1312. <https://doi.org/10.1016/j.ard.2025.04.011>

Important note

To cite this publication, please use the final published version (if applicable).
Please check the document version above.

Copyright

Other than for strictly personal use, it is not permitted to download, forward or distribute the text or part of it, without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license such as Creative Commons.

Takedown policy

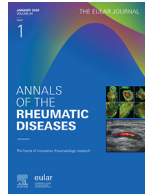
Please contact us and provide details if you believe this document breaches copyrights.
We will remove access to the work immediately and investigate your claim.



ELSEVIER

Contents lists available at ScienceDirect

Annals of the Rheumatic Diseases

journal homepage: <https://www.sciencedirect.com/journal/annals-of-the-rheumatic-diseases>

Rheumatoid arthritis

Time-independent disease state identification defines distinct trajectories determined by localised vs systemic inflammation in patients with early rheumatoid arthritis

Nils Steinz^{1,*}, Tjardo D. Maarseveen¹, Erik B. van den Akker^{2,3}, Andrew P. Cope⁴, John D. Isaacs⁵, Aaron R. Winkler⁶, Tom W. J. Huizinga¹, Yann Abraham⁷, Rachel Knevel^{1,3,8}

¹ Department of Rheumatology, Leiden University Medical Centre, Leiden, the Netherlands

² Leiden Computational Biology Centre, Leiden University Medical Center, Leiden, the Netherlands

³ The Delft Bioinformatics Lab, Pattern Recognition and Bioinformatics, Delft University of Technology, Delft, the Netherlands

⁴ Centre for Rheumatic Diseases, School of Immunology and Microbial Sciences, Faculty of Life Sciences and Medicine, King's College London, London, UK

⁵ Musculoskeletal Unit, Newcastle Upon Tyne Hospitals NHS Foundation Trust, Newcastle upon Tyne, UK

⁶ Inflammation and Immunology Research Unit, Pfizer Inc, Cambridge, MA, USA

⁷ Janssen Research and Development, Beerse, Belgium

⁸ Translational and Clinical Research Institute, Newcastle University, Newcastle upon Tyne, UK

ARTICLE INFO

Article history:

Received 8 November 2024

Received in revised form 14 March 2025

Accepted 8 April 2025

ABSTRACT

Objectives: Patients with rheumatoid arthritis (RA) display different trajectories towards improvement of disease. We aimed to disentangle the heterogeneity of RA disease trajectories from the first clinical visit onwards using graph-based pseudotime analysis.

Methods: We studied early patients with RA over 1.5 years in 2 data sets: Leiden (Netherlands), n = 1237, with 5017 visits, and Towards a Cure for Early Rheumatoid Arthritis (TACERA) (United Kingdom), n = 243, with 750 visits. We created a pipeline for time-independent clustering of clinical and haematologic features to identify disease states. Sequence analyses of these states defined the trajectories. We studied the predictability of the trajectories with baseline features.

Results: Clustering identified 8 disease states with localised inflammation (joints) and systemic inflammation (erythrocyte sedimentation rate [ESR] or leucocytes) as the main discriminating factors. The disease state sequences consisted of 4 trajectories, which we independently replicated in TACERA: A, high ESR; B, rapid progression from many inflamed joints towards remission; C, high leucocytes; and D, many inflamed joints with poor prognosis. Systemic vs local inflammation patterns showed moderate predictability at baseline (sensitivity of 71% and precision of 0.73 for trajectory A, although lower precision of 0.52 for trajectory B), while other trajectories were less predictable. Trajectories C and D had strong resemblance with B at baseline but deteriorated into less favourable trajectories. Patients in trajectory A were more often female and on average older. The trajectories were not explained by time till disease-modifying anti-rheumatic drug, baseline disease activity, or symptom duration. The suboptimal trajectories coincided with worse patient-reported outcomes, even when the inflammation was mainly systemic.

*Correspondence to Nils Steinz.

E-mail address: n.stein@lumc.nl (N. Steinz).

Handling editor Josef S. Smolen.

<https://doi.org/10.1016/j.ard.2025.04.011>

Conclusions: We identified 4 distinct trajectories in early RA, differentiating RA into localised vs systemic inflammation. Our results highlight potential differences in disease pathology and opportunities for further targeted treatment. Inevitably, patterns without linkage to our selected features could not be detected.

WHAT IS ALREADY KNOWN ON THIS TOPIC

- The course of early rheumatoid arthritis (RA) after disease-modifying antirheumatic drug (DMARD) initiation differs in slow, moderate, and rapid progressors.
- Longer symptoms duration before DMARD initiation, older age, female sex, and seropositivity are associated with worse outcomes.

WHAT THIS STUDY ADDS

- Early RA consists of at least 4 different trajectories characterised by systemic inflammation (elevated erythrocyte sedimentation rate [ESR] or leucocytes) and localised inflammation (either fast or no/bad responders to treatment).
- Patients with mainly systemic inflammation (elevated ESR or leucocytes) are less likely to reach remission and have worse patient-reported outcomes.
- Most patients' trajectories are predictable with baseline joint and blood measurements and are not further explained by symptom duration, baseline disease activity, or timing of methotrexate prescription.

HOW THIS STUDY MIGHT AFFECT RESEARCH, PRACTICE OR POLICY

- To achieve optimal disease control, it is important to suppress inflammation in both the joints and systemic circulation.
- The trajectories underline the relevance of studies to the pathological differences of systemic vs local inflammation.
- Our study highlights a subgroup who deviates from the favourable trajectory who might benefit from more intensive treatment.

INTRODUCTION

Despite advances in RA treatment, disease trajectories from first clinical visit to remission vary significantly [1–5]. Furthermore, the path to remission is rarely linear, with patients frequently experiencing fluctuations in disease activity over time [5]. To identify optimal moments for intervention and improve overall treatment outcomes, it is crucial to capture and analyse real-life trajectories of patients with early-treated RA.

Previous research identified 3 trajectories: rapid, gradual, and slow progressors [6,7]. These studies categorised patients exhibiting comparable disease activity difference at equivalent postbaseline time points [8]. By design, these studies could not have detected the granularity of disease trajectories where patients can move back to previous disease states. To further improve the knowledge on disease trajectory, we need to capture the granularity, apply methods that do not assume linear or polynomial relationships, and allow for variable transition speeds through disease states.

Unsupervised clustering methods are able to uncover unknown or invisible structures. They use functions and mathematical equations to separate data points, rather than using labels to differentiate data. These equations establish either linear or nonlinear correlations between data points, thereby forming clusters primarily based on the similarity between variables [9–11]. Specifically, graph-based pseudotime analysis has proven its

value to discern the order of cell states in single-cell analyses [12]. Clinical patient data are different from single-cell data as there are no clear markers, highlighting whether a state is early or late, with the potential exception of the remission state. Certain states (eg, modest amount of joint inflammation) can occur either early or late during the disease course. Luckily, with the accessibility of electronic health records and observational cohorts, longitudinal data are available where the sequence of disease states can be deduced from the chronology of clinical visits.

In this study, we aimed to capture the real-world disease trajectories of patients with RA during their first 1.5 year at the outpatient clinic. First, we identified existing disease states independent of time integrating swollen joint counts (SJC) and tender joint counts (TJC) with multiple standard laboratory measurements. Subsequently, we defined disease trajectories by grouping patients with similar disease state sequences. The ultimate goal of this research was to elucidate distinct, explainable patient trajectories that elucidate the heterogeneity of early RA.

METHODS

Data

Our retrospective study contained 2 independent longitudinal datasets: the Dutch Leiden electronic health records (EHR) and the Towards a Cure for Early Rheumatoid Arthritis (TACERA) cohort from the United Kingdom. Both data sets consist of new/recent onset adult patients with RA (≥ 18 years), disease-modifying antirheumatic drug (DMARD) naive at baseline.

The Leiden EHR cohort contains patients who were newly referred to the Leiden outpatient clinic between 2011 and 2022 and received a diagnosis of RA within the first 1.5 years of follow-up by their rheumatologist [13]. Our data contains the first 1.5 years of records from patients who had at least 2 visits, including baseline. A single visit was defined as a haematologic laboratory test and a consultation with a rheumatologist within a 10-day window and incomplete visits were disregarded. A total of 1237 patients met the inclusion criteria for the study and collectively had 5017 visits. The treatment of patients reflects clinical care.

The Towards A Cure for Early Rheumatoid Arthritis (TACERA) cohort consists of newly diagnosed patients with RA (European Alliance of Associations for Rheumatology (EULAR) 1987 or American College of Rheumatology (ACR) 2010) [14,15] with a maximum of 12-month symptom duration who were DMARD naive and seropositive (243 patients with 750 total visits) [7]. They were recruited in 2014 at time of DMARD initiation and the follow-up ranged from 0.5 to 1.5 years. Treatment changes were in accordance with the National Institute for health and Care Excellence guidelines [16].

Feature selection for clustering

Feature selection was guided by 3 key principles: (1) focus on direct measures of disease activity that can change over time, (2) availability of consistent measurements across routine clinical visits, and (3) exclusion of treatments and their related variables and measurements. This led us to include SJC28 and TJC28,

erythrocyte sedimentation rate (ESR), leucocyte count, haemoglobin (Hb), mean corpuscular volume (MCV), and thrombocyte counts as measures of systemic inflammation. These variables capture the core clinical presentation of disease stages of RA and represent features routinely used in clinical decision-making. Anti-cyclic citrullinated peptide (aCCP) and rheumatoid factor (RF) were not part of the clustering as they did not meet the principles. In accordance with the third principle, the clustering excluded medication, dosage, and liver enzymes.

Clustering of independent visits to identify disease states

Figure 1 provides a schematic representation of the study’s pipeline from raw data to identification of disease states and disease trajectories. Visits were excluded if they missed clustering variables, occurred beyond 1.5 years after initial rheumatologist visit, or when leucocytes exceeded 25×10^9 (as this was most likely not RA related). Patients were excluded if they lacked complete baseline data or had fewer than 2 visits within 1.5 years of follow-up. The input data were the SJC28, TJC28, ESR, leucocyte count, Hb, MCV, and thrombocyte counts at each visit. C-reactive protein was not included due to low availability, and for the TACERA, the MCV was excluded due to its absence in the data set. The features were normalised using minimum to maximum normalisation between 0 and 1. For identification of the disease states, we clustered the visits using the Phenograph [17] with the cosine similarity metric and the Leiden algorithm. The time of the visits, sex, age, and seropositivity were not used for clustering.

We determined the number of clusters (disease states) based on data clustering consistency and cluster size. For this, we set the optimal k for the k -nearest neighbours algorithm in the Phenograph package by testing a range of k values between 10 and 400, with increments of 10. We compared the different configurations for the clustering model to find the most stable region in the adjusted Rand index between k values. For transferability testing, we used machine learning algorithms (logistic regression, support vector machine, and random forest) to classify the clusters based on the features.

Defining patient trajectories

With the real timing of visits, we could order the disease states for each patient creating individual trajectories. We identified commonalities in patients’ trajectories with the adjusted Bobroske algorithm [18] to disallow transposition. Matches were calculated as the inverse logarithm of the frequency of the disease state occurrence and further normalised to within a range between 1 and 5. Mismatches were transformed by multiplying the matches by -0.5 , and indels were calculated as mismatches multiplied by 0.75 . The sequence score was transformed using the same calculation as in the original article. Pairwise sequence similarity was applied to all patients, resulting in a similarity matrix. The Phenograph package’s clustering algorithm was then applied to cluster similar trajectories.

After identifying trajectories, we performed post hoc analyses to characterise these patterns using variables excluded from the clustering analysis. This included examining medication usage, treatment survival, and patient characteristics (eg, age, CCP, and RF) across the different trajectories.

Clinical differences between trajectories

To compare differences in patient-reported outcome measure (PROM), we created a linear mixed effect model with random intercept. The model was fitted with fixed-effects age, sex, months, and trajectory cluster with an interaction between months and trajectory cluster. We tested whether the interactions between month and trajectory clusters were significant.

Baseline features used to predict the future trajectory

In order to predict the trajectory that a patient would take, we trained an support vector machine (SVM) model using crossvalidation on 80% of the data and used a 20% test set for final results. The model used the baseline features that were used to cluster the data. We then analysed the discrepancies between the baseline trajectory assignment and their true trajectory. To investigate the impact of age, sex, seropositivity, days until medication, and symptom duration on the

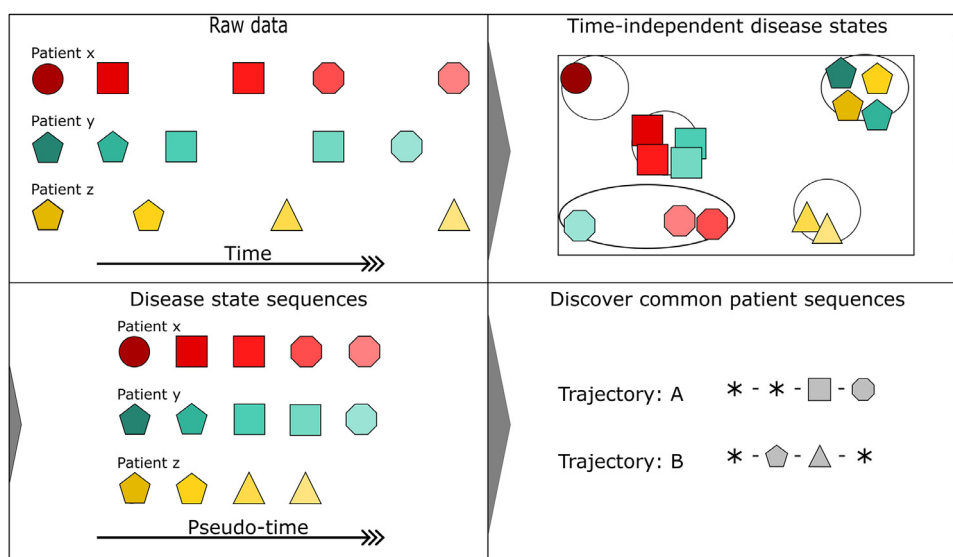


Figure 1. Individual patient visits are clustered based on similar phenotypic characteristics, irrespective of patient, treatment, or time. Next, the sequence of patient states is inferred from the date of visits, leading to individual patient trajectories. Disease trajectories are defined based on a shared sequence of states across patients using sequence similarity and clustering. The colours indicate different patients, and the hue decreases with time in a patient visit sequence. Different shapes indicate different disease states.*Any cluster could be in that time point.

prediction of trajectories, we used logistic regression on the subset of patients with known symptom duration.

Treatment survival analysis

To assess whether differences in outcomes between trajectories reflect patterns of treatment response, we performed survival analyses of DMARD persistence. We performed 2 analyses: first DMARD survival and second DMARD survival. For the first DMARD analysis, time zero was set at DMARD initiation, and events were defined as switching to or adding another DMARD. For the second DMARD analysis, we calculated survival until the start of the third DMARD change event if patients had a second DMARD. In both analyses, patients were censored at the last follow-up visit if no switch occurred. We used Kaplan–Meier to estimate survival probabilities and compared trajectories using log-rank tests.

RESULTS

Baseline description and general patient trajectory

The 2 data sets had some differences due to inclusion criteria (Table 1). The Leiden patients were slightly older and less often female. As per study design, patients in TACERA were seropositive. The overall seropositivity rate was 55.2% in Leiden. At the first visit, the majority of patients in Leiden (53%) received a DMARD, which rapidly increased over the follow-up visits (Supplementary Fig S1). In TACERA, this was 100%, by design. On average, both the Leiden and TACERA populations improved over time on all features (Table 1). Patient characteristics of those excluded from the Leiden cohort due to missing data are presented and compared in Supplementary Table S1.

Patient disease states during the first 1.5 years of RA

We used a data-driven approach as outlined earlier to define the optimal number of disease states (clusters). This clearly showed 8 clusters over any other number of clusters

(Supplementary Fig S2). We then further confirmed the stability of the clusters with an SVM algorithm, which predicted the clusters with 94.8% accuracy on a holdout test set (Supplementary Fig S3).

The 8 identified clusters described the disease states exhibited by patients during their visits to the clinic (Fig 2A and Supplementary Table S2). Although time of visit was not part of the clustering, the disease states exhibited significant temporal differences (Fig 2B). The disease state that contained the most late-time visits was L1. At this disease state, patients had the lowest disease activity with near-zero number of swollen and tender joints, lowest ESR, and within normal range leucocytes and Hb. L8 is predominantly an early visit state that contains no visits after 1 year. At this disease state, patients have low swelling and tenderness and low Hb.

For the in-between disease states, the main differentiation is based on the amount of swollen and tender joints vs systemic inflammation (ESR or leucocytes).

Among the high joint inflammation states, L7 is the most severe state. Patients in this state have a high number of swollen and tender joints and a moderately elevated ESR. The L7 disease state not only is frequently observed at baseline but also occurred less frequently at later stages of disease course. A similar but more moderate disease state L4 is characterised by numerous swollen and tender joints and without elevated ESR. L4 is common early in the disease course and common after 1 year of follow-up.

The inflamed disease states could be further differentiated into elevated ESR or leucocytes. The elevated ESR states were L6, L5, and L2. L6 contains patients who have moderately swollen and tender joints and a highly elevated ESR. Patients with L6 disease were commonly seen at baseline, and some may present with the disease state later on in their treatment. Patients in disease state L5 have similar disease characteristics, visit frequency, and time after baseline to L6, but with a less elevated ESR. L2 contains patients at a low disease state: low to no swollen and/or tender joints but still a modestly elevated ESR. This elevated ESR coincides with a subtle increase in leucocytes, thrombocytes, and decreased Hb. This condition could be

Table 1
General data description of the Leiden cohort and the TACERA cohort

Features	Leiden		TACERA		P	
	Baseline	Last visit	Baseline	Last visit	Baseline	Last visit
No. of patients	1237		243			
No. of visits	4 (2-5)		3 (2-4)			
Female	797 (64.4)		177 (72.8)		<.05	
RF positive	570 (46.1)		207 (85.2)		<.05	
aCCP positive	571 (46.2)		223 (91.8)		<.05	
Seropositive	683 (55.2)		243 (100)		<.05	
Age (y), mean (SD)	59.8 (14.7)		53.0 (14.9)		<.05	
TJC28 ^a	4 (2-8)	1 (0-3)	9 (4-14)	2 (0-5)	<.05	<.05
SJC28 ^a	3 (1-7)	0 (0-2)	6 (3-10)	1 (0-3)	<.05	.786
ESR ^a	28 (11-45)	14 (6-29)	27 (12-43)	12 (6-24)	.228	<.05
Leucocytes ^a	7.9 (6.3-9.6)	7.2 (6.0-8.8)	7.8 (6.6-9.2)	6.3 (5.1-7.8)	.299	<.05
Haemoglobin ^a	8.2 (7.5-8.8)	8.3 (7.7-8.8)	8.2 (7.5-8.8)	8.2 (7.6-8.8)	<.05	.221
MCV ^a	90 (87-93)	92 (89-96)	NA	NA	NA	NA
Thrombocytes ^a	290 (246-352)	264 (225-314)	293 (255-365)	267 (226-312)	.199	.885
DAS28	4.3 (3.3-5.3)	2.9 (2.1-3.9)	5.0 (4.1-5.8)	3.0 (2.0-4.0)	<.05	.651
DMARD assignment	650 (52.5)	1148 (92.8)	243 (100)	243 (100)	<.05	<.05

aCCP, anticyclic citrullinated peptide; DAS, disease activity score; DMARD, disease-modifying antirheumatic drug; ESR, erythrocyte sedimentation rate; MCV, mean corpuscular volume; RF, rheumatoid factor; SJC, swelling joint count; TACERA, Towards a Cure for Early Rheumatoid Arthritis; TJC, tender joint count.

Values are median (IQR) or n (%) unless specified. Welch t test was applied for the P values.

^a Features included for clustering.

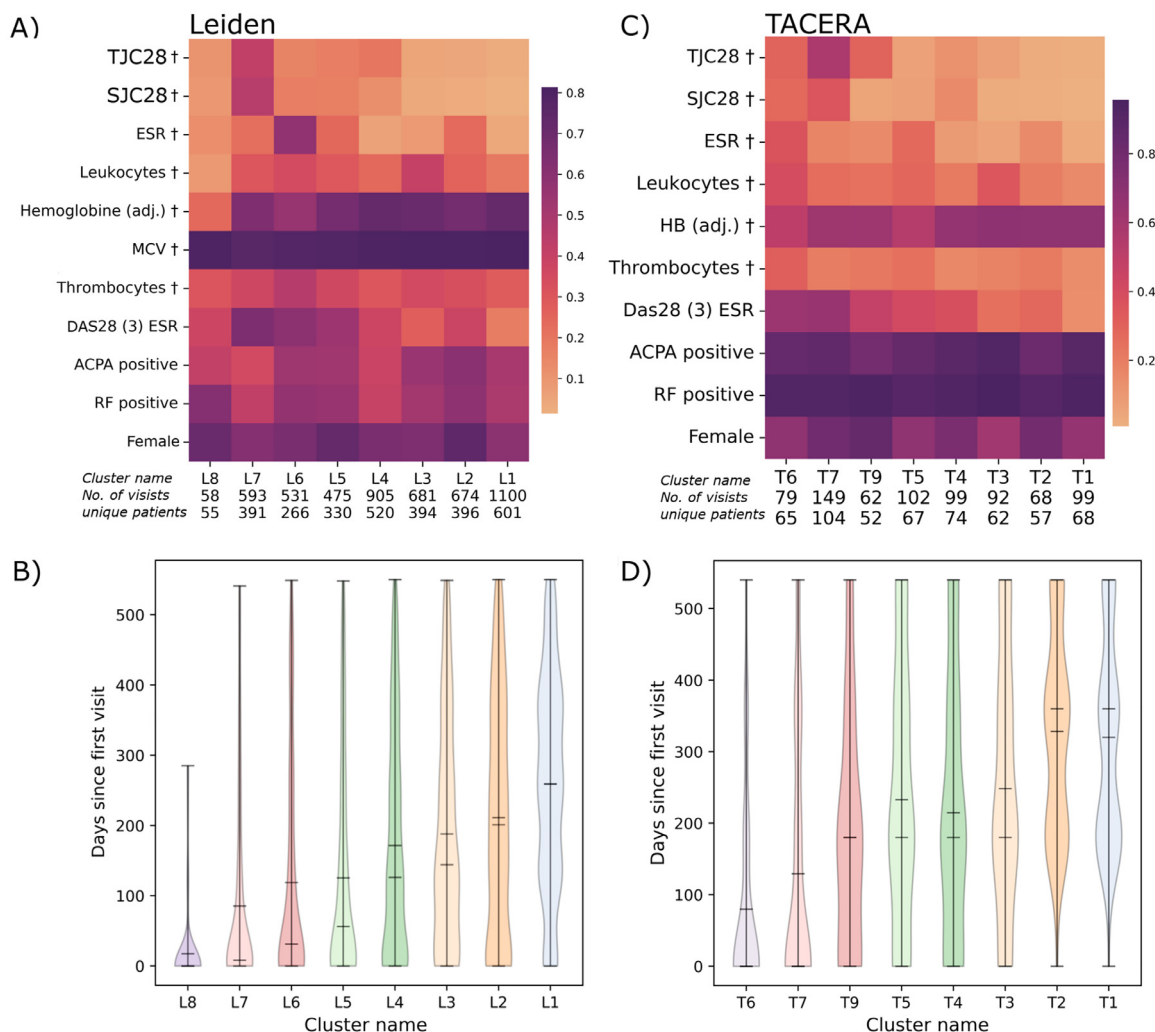


Figure 2. Results of the clustering and sequence analyses on the Leiden or TACERA data. (A,C) Heatmap of mean characteristics of the disease states (clusters). All values are scaled between 0 and 1 (orange–purple). (B,D) Number of days the visit took place after baseline for each cluster. †Variables used for clustering. TACERA, Towards a Cure for Early Rheumatoid Arthritis.

observed in the clinic in all states within the first 1.5 years, but the majority of the L2 assigned visits were after 200 days of first visit.

L3 describes a phenotype with a leucocyte count above the upper reference values (mean, 11.0×10^9 ; SD, 2.69; reference = $4.0\text{--}10.0 \times 10^9$) and a low amount of joint and blood inflammation. This state of the disease was most prevalent 1 year after the baseline.

Similar disease states observed in independent replication in TACERA

For the TACERA data set, we applied the same normalisation and unsupervised clustering pipeline as used for Leiden. To facilitate comparability of cluster content with Leiden, we set the *k* to find 8 clusters in TACAERA. The content of the clusters were defined in an unsupervised manner, without using any knowledge from the findings in the Leiden data.

Within TACERA, we saw that the disease states resembled those from the Leiden data (Fig 2C). The optimal disease state L1 was also present in TACERA which we named T1. T1 shows no active disease with a near-zero amount of swollen and tender joints and no elevation in ESR.

All first 7 disease state (L1–L7) were also present in TACERA. We identified a distinct state characterised by patients with tender joints and mildly elevated ESR (<30) in TACERA. While this

state shares some features with L4 (tender joints), it represents a unique intermediate phenotype not observed in the Leiden cohort. To avoid confusion with L8 (the early disease state unique to Leiden), we designated this TACERA-specific state as T9. Similar as in the Leiden data, the main differentiating factor in TACERA is based on local (joints) vs systemic inflammation (ESR and leucocytes).

Different disease trajectories in patients with early RA

At the group level, we observed a clear pattern of reduction of disease activity towards reaching remission; 47.0% and 28.0% of the patients reached the optimal disease state L1/T1 within 1.5 years, and 31.4% and 23.5% got close to the optimal state and ended in L2/T2. This suboptimal end state appears optimal based on SJCs and TJCs, but the ESR was still elevated. This deviation from the optimal disease state was rather subtle as most values were within the normal range.

To find subsets of trajectories, we independently clustered the patient-level sequences of Leiden and TACERA and found 4 separate trajectories in Leiden and 5 in TACERA (Fig 3 and Supplementary Fig S4). These revealed a more granular pattern of disease improvement and flare than the average patterns. At baseline, patients across the trajectories differed in age, sex, seropositivity, baseline disease activity score (DAS)28, and symptom duration (Tables 2 and 3).

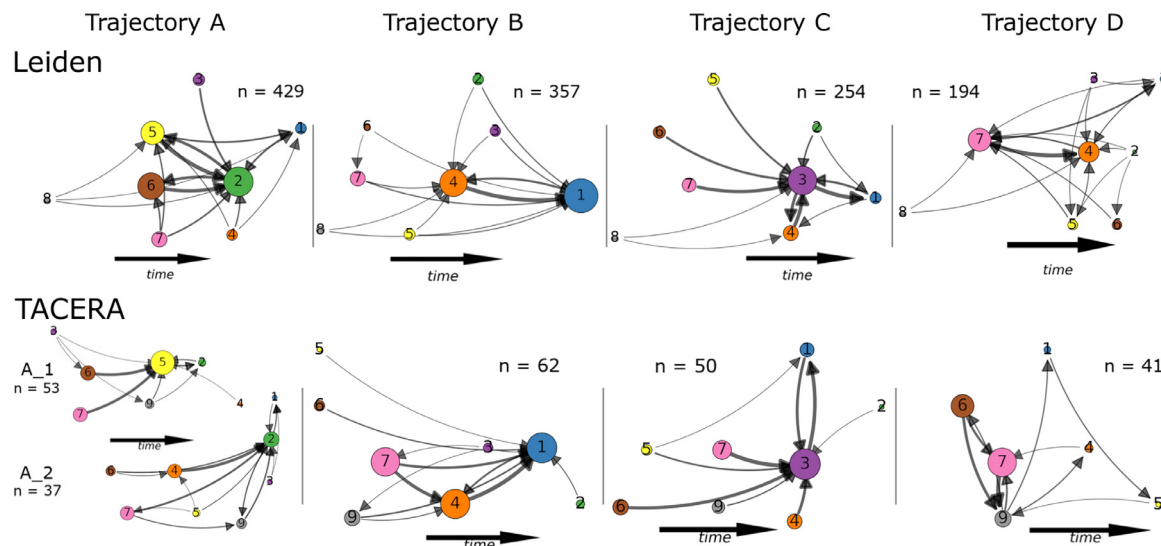


Figure 3. Trajectories based on the cluster sequence per patient. Four common trajectories were identified and were distinguished by (A) erythrocyte sedimentation rate, (B) good trajectory, (C) leucocytes, and (D) bad trajectory. The nodes are the disease states with the size of the node marking the number of patient visits. The edge thickness reflects the transition frequency. The nodes on the left side of the graph are, on average, more prevalent at an earlier state of the trajectory, while the nodes on the right side are more common at a later state. TACERA, Towards a Cure for Early Rheumatoid Arthritis.

Although independently analysed, the trajectories from Leiden and TACERA exhibited strong similarities as summarised further.

Leiden-A, TACERA-A.1, and TACERA-A.2—high ESR clusters at start and end (L6 to L5 to L2, and T6 to T5 to T2)

While patients following this trajectory could start in any disease state, it was most common for them to start in the active disease state L/T6 or L/T5 with elevated ESR and active swollen and tender joints and end up in the suboptimal disease state L/T2. A small subset of patients reached the optimal final state (L/T1), but most moved back to L/T2. Within TACERA, trajectory A was subdivided into 2 similar trajectories marked by elevated

ESR, where patients in TACERA-A.1 have the worse outcome by remaining in the modestly active disease T5, instead of moving towards T2. Trajectory A was predominantly female. The prevalence of CCP and RF was higher than average in Leiden and more evenly distributed in TACERA.

Leiden-B and TACERA-B—rapid progression from primarily joint inflammation towards the ideal cluster (L1 and T1)

Patients can start from any cluster and reached L/T1 within a year, either directly or via L/T4. Trajectory B had the most marked DAS28 improvement. They were more frequently male were younger. On average, it took patients in trajectory B

Table 2
Descriptive characteristics of Leiden patients and trajectories

Feature	Leiden patients	Leiden trajectory A	Leiden trajectory B	Leiden trajectory C	Leiden trajectory D
No. of patients	1237	429	357	254	197
Visits	4 (2-5)	3 (2-5)	4 (3-5)	4 (3-5)	3 (2-5)
Female	797 (64.4)	309 (72.0)	214 (59.9)	155 (61.0)	119 (60.4)
Age (y), mean (SD)	59.8 (14.7)	63.01 (13.8)	56.03 (14.27)	58.76 (16.45)	61.1 (13.45)
Rheumatoid factor	570 (46.1)	217 (50.6)	162 (45.4)	123 (48.4)	68 (34.5)
aCCP	571 (46.2)	227 (52.9)	160 (44.8)	130 (51.2)	54 (27.4)
Seropositive	683 (55.2)	263 (61.3)	195 (54.6)	145 (57.1)	80 (40.6)
Baseline TJC28 ^a	4 (2-8)	4 (2-7)	3 (1-6)	3 (1-7)	9 (5-15)
Baseline SJC28 ^a	3 (1-7)	3 (1-6)	2 (0-5)	2 (1-6)	7 (4-13)
Baseline ESR ^a	28 (11-45)	43 (31-67)	14 (6-25)	19 (11-36)	19 (9-36)
Baseline DAS28(3) ESR, mean (SD)	4.27 (1.38)	4.71 (1.12)	3.61 (1.3)	3.94 (1.4)	4.89 (1.39)
Last visit TJC28 ^a	1 (0-3)	2 (0-4)	1 (0-2)	1 (0-2)	4 (1-7)
Last visit SJC28 ^a	0 (0-2)	1 (0-3)	0 (0-1)	1 (0-2)	3 (0-5)
Last visit ESR ^a	14 (6-29)	32 (19-45)	6 (2-11)	11 (6-19)	9 (6-17)
Last visit DAS28(3) ESR, mean (SD)	3.01 (1.29)	3.72 (1.15)	2.13 (0.95)	2.76 (1.15)	3.37 (1.23)
Symptom duration at baseline (d)	155 (62-548)	154 (62-730.0)	180 (62-546)	233 (80-730)	115 (42-276)
Time till first DMARD (d)	28 (14-105)	28 (14-120)	28 (14-101)	28 (14-147)	28 (14-65)
Prednisolone used ever (%)	72.8	77.2	65.5	74.6	78.3
Most prevalent first DMARD (%)	MTX (80.0)	MTX (81.0)	MTX (82.5)	MTX (80.6)	MTX (78.6)
Remission ^b (%)	56.9	25.4	90.5	72.4	44.7

aCCP, anticyclic citrullinated peptide; DAS(3), 3-component disease activity score for the 28 joint scheme; DMARD, disease-modifying antirheumatic drug; ESR, erythrocyte sedimentation rate; MTX, methotrexate; SJC, swollen joint count; TJC, tender joint count.

Values are median (IQR) or n (%) unless specified.

^a Features included for clustering.

^b Reached DAS28(3) ESR < 2.6 at some point in the trajectory.

208 days to reach low disease activity or remission, while starting with an average of 3 tender and swollen joints at baseline.

Leiden-C and TACERA-C—transition through or ending at high leucocyte state (L3 and T3)

Patients in trajectory C experience difficulty reaching cluster L1 and may divert from L3 to L4, which can be considered a worsening disease activity. Time between symptom onset and the first visit was the longest for this group, while other aspects of the data set were average. The main difference between trajectory C and the optimal trajectory B was an increase in leucocytes. Patients in trajectory C had the highest leucocyte count before and after prednisone (Supplementary Fig S5).

Leiden-D and TACERA-D—poor prognosis (L7 to L4 and T7 to T4)

Trajectory D was the most severe trajectory, as patients do not reach the optimal disease state L/T1 or even the low disease states L/T2 and L/T3. Instead, they end up in an active disease state of L/T4 with inflamed joints but low levels of blood inflammation. Trajectory D had the worst baseline TJC and SJC. The DAS decreased from 4.65 to 3.54 in Leiden and 5.75 to 4.60 in TACERA, indicating a moderate disease activity state, over the course of 1.5 years. The patients in this trajectory exhibited the shortest symptom duration until their first visit to the clinic, with an average of 101 days. Additionally, trajectory D exhibited the lowest prevalence of seropositive patients. At the start, patients in trajectory D were similar to those in trajectory B.

Most patients in the trajectories received DMARDs within the first month, with a median time to DMARD of 15 (14–19) days. In all trajectories, methotrexate (MTX) was the primary DMARD, with only slight variations in average usage. In addition, compared with the other trajectories, patients in trajectory B receive least frequently prednisolone. A comparable proportion of patients within the TACERA cohort receive MTX at baseline, with the remainder receiving other conventional synthetic

DMARDs. However, TACERA patients were less likely than Leiden patients to receive steroids (triamcinolone or prednisone), with the highest percentage being trajectory D.

Patient-reported outcomes differ between trajectories

The trajectories differed in both the disease states patients' transition through as well as the final disease state. The question arises whether the differences between the trajectories in systemic or localised inflammation are clinically relevant, particularly when the difference was driven by ESR or leucocytes and not per se the joints. To analyse the differences in patient-perceived progressions and correlation with trajectories, we examined the PROMs from the TACERA cohort, which were available for 72.4% of the visits. Our findings showed that the trajectories that do not lead to the optimal state (L/T1) also resulted in worse outcomes in the PROMs (Fig 4). The disease progression and activity described by DAS28 (A) and simplified disease activity index (SDAI) (B) showed similar trajectories with the SDAI being worse overall. In addition, in PROMs, trajectory B was the most favourable trajectory. Here, patients start with clear impairment and steadily improve over time as measured in health assessment questionnaire (HAQ), fatigue, and global assessment. Trajectory A_1 and A_2 (the trajectories characterised by increased ESR) have fairly good improvement in HAQ, but their fatigue remains unchanged. Trajectory D, the poorest trajectory with highly inflamed joints, starts with the poorest HAQ, slightly improving in the first year, but returns to HAQ levels similar to baseline. This trajectory differs significantly from B (Supplementary Table S4). For the fatigue score, the plot lines were stable, with only groups B and C demonstrating a distinct downwards trend in the first 6 months. Surprisingly, trajectory C returns to the initial values, while B continues to improve. For both patient global assessment (PGA) and evaluator global assessment (EGA), trajectory C showed a significantly different

Table 3
Descriptive characteristics of TACERA patients and trajectories

Feature	TACERA patients	TACERA trajectory A_1	TACERA trajectory A_2	TACERA trajectory B	TACERA trajectory C	TACERA trajectory D
No. of patients	243	53	37	62	50	41
Visits	3 (2-4)	3 (2-4)	3 (2-4)	4 (3-4)	3 (2-4)	3 (2-4)
Female	177 (73.5)	36 (67.9)	30 (81.1)	45 (72.6)	32 (64.0)	34 (82.9)
Age (y), mean (SD)	53.4 (15.00)	60.23 (12.4)	53.41 (13.77)	49.58 (15.00)	50.48 (16.23)	51.24 (14.56)
Rheumatoid factor	223 (91.8)	47 (88.7)	34 (91.9)	58 (93.5)	48 (96.0)	36 (87.8)
aCCP	207 (85.5)	43 (81.1)	30 (81.1)	53 (85.5)	45 (90.0)	37 (90.2)
Seropositive	100	100	100	100	100	100
Baseline TJC28 ^a	9 (4-14)	7 (3-13)	8 (4-15)	9 (5.2-14)	6.5 (3-10)	13 (9-17)
Baseline SJC28 ^a	6 (3-10)	6 (3-10)	7 (5-9)	6 (4-9.8)	4 (2-7)	8 (4-10)
Baseline ESR ^a	27 (12-43)	40 (27-52)	26 (14-44)	17 (8-32)	15 (7-31)	30 (17-65)
Baseline DAS28(3) ESR, mean (SD)	5.18 (1.37)	5.2 (1.12)	5.02 (1.38)	5.11 (1.32)	4.17 (1.27)	5.75 (1)
Last visit TJC28 ^a	2 (0-5)	1 (0-3)	0 (0-2)	1 (0-3)	1 (0-3.8)	8 (6-13)
Last visit SJC28 ^a	1 (0-3)	1 (0-2)	0 (0-2)	0 (0-2)	0 (0-1.8)	4 (2-6)
Last visit ESR ^a	12 (6-24)	27 (17-33)	21 (10-24)	7 (2-11)	6 (2.2-9.8)	16 (8-25)
Last visit DAS28(3) ESR, mean (SD)	3.04 (1.36)	3.36 (0.91)	2.93 (0.99)	2.40 (1.33)	2.30 (1.10)	4.60 (1.00)
Symptom duration at baseline (d)	137 (87-198)	119 (84-159)	138 (93-182)	146 (92-199)	153 (93-218)	173 (83-259)
Time till first DMARD (d)	0 (0-0)	0 (0-0)	0 (0-0)	0 (0-0)	0 (0-0)	0 (0-0)
Steroid used ever (%)	72.4	81.1	64.9	72.6	68.0	73.2
Remission ^b (%)	38.4	20.8	32.4	62.9	62	0

aCCP, anticyclic citrullinated peptide; DAS(3), 3-component disease activity score for the 28 joint scheme; DMARD, disease-modifying antirheumatic drug; ESR, erythrocyte sedimentation rate; SJC, swollen joint count; TACERA, Towards a Cure for Early Rheumatoid Arthritis; TJC is tender joint count.

Values are median (IQR) or n (%) unless specified.

^a Features included for clustering.

^b Reached DAS28(3) ESR < 2.6 at some point in the trajectory.

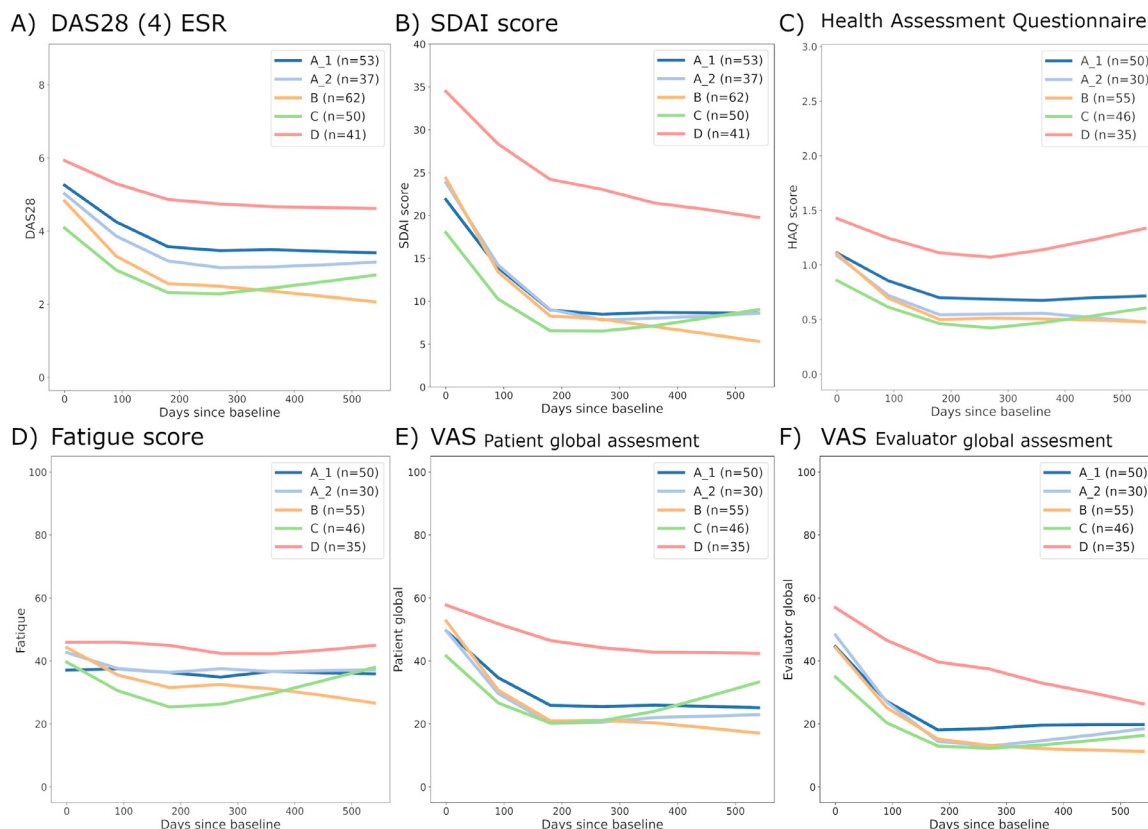


Figure 4. Evolution of patient-reported outcomes over time. Trajectories A_1 and A_2 are characterised by elevated erythrocyte sedimentation rate (ESR) levels, with A_1 associated with an unfavourable prognosis. Trajectory B represents the optimal outcome, while Trajectory C is categorised by elevated leucocyte counts. Trajectory D represents the least favourable progression group. The graphs depict (A) disease activity score (DAS)28 ESR trend, (B) simplified disease activity index (SDAI) trend, (C) the trend of health assessment questionnaire score, (D) fatigue reported by patients, (E) disease activity experienced by patients; (F) disease activity assessed by the rheumatologist. VAS, visual analogue scale.

path after 1 year compared with B. For trajectory D, only PGA was significantly different.

Predicting trajectories based on baseline features

We built a classifier to predict patients' trajectory using baseline value of the features that were used for defining the disease states. This model exhibited an overall accuracy of 60%. Particularly trajectories A and B could be accurately predicted with a sensitivity (recall) of 71% for both (Fig 5 and Supplementary Table S3), while trajectory B was less precise with a precision of 0.52 vs 0.73 for trajectory A. The prediction of D and, to a lesser extent C, at baseline had lower sensitivity rates (26% and 48%, respectively), with D additionally not being precise (0.38) compared with C (0.71). Patients who ultimately end up in trajectory C or D were frequently labelled or predicted to follow trajectory B or even A. Many patients at baseline exhibited characteristics similar to those of patients following trajectory B. Yet, they ultimately deviate towards the less favourable trajectories C and D. Similarly, approximately 15% of patients in different trajectories have a similar baseline compared with those in trajectory A. Women, older people, and people with anti-CCP (aCCP) were more likely to end up in trajectory A than B (Supplementary Fig S6). Age, sex, and aCCP improved the predictive model, while MTX, time to DMARD, and symptom duration until the first visit did not significantly improved the prediction of any of the trajectories at baseline (Supplementary Tables S5 and S6).

Treatment persistence across trajectories

To understand treatment patterns across trajectories, we analysed both first and second DMARD survival (Fig 6). First DMARD survival at 1.5 years varied significantly between trajectories, with trajectory B showing the highest persistence (62% remaining on initial DMARD) compared to trajectories A (54%; $P = .002$), C (58%; $P = .924$), and D (47%; $P = .001$). Notably, while trajectories B and C showed similar persistence ($P = .924$), both demonstrated significantly better drug survival than trajectory D (C vs D: $P = .002$) and trajectory A (C vs A: $P = .008$). Trajectories A and D showed comparable persistence rates ($P = .278$).

This pattern was largely maintained for second DMARD survival, where trajectory B again showed superior persistence (68%) compared with trajectories A (58%; $P < .001$), C (60%; $P = .559$), and D (51%; $P = .003$). Trajectory C maintained significantly better persistence than both trajectories A ($P = .009$) and D ($P = .030$), while trajectories A and D showed similar survival rates ($P = .920$).

DISCUSSION

Our findings confirm the existence of distinct disease trajectories in early RA. We uncovered 4 RA trajectories: (A) elevated ESR, (B) localised joint inflammation with favourable outcome, (C) elevated leucocyte count, and (D) persistent localised joint inflammation. These trajectories can be broadly categorised into localised inflammation and systemic inflammation patterns,

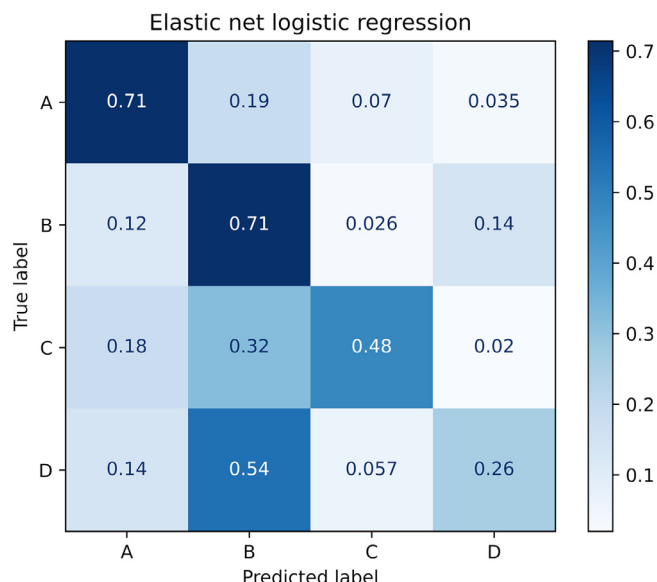


Figure 5. The confusion matrix on the test set illustrates the predictive accuracy of the trajectory based on the features used in the clustering process at baseline. The predictions for A and B are generally accurate, while those for C are less accurate and those for D are not predictable based on the baseline characteristics and features (tender joint count 28, swollen joint count 28, erythrocyte sedimentation rate, and other laboratory values).

with systemic inflammation generally associated with poorer outcomes. The study’s strengths include a reasonable sample size of patients with DMARD-naïve early RA, independent replication in a second data set, additional validation with PROMs, and the use of widely available model parameters.

Replication was crucial to validate our unsupervised clustering and avoid detecting spurious patterns. Our replication data, TACERA, differed by containing only seropositive patients fulfilling RA classification criteria at baseline, resulting in a more severe RA population. Applying the unsupervised clustering to

TACERA allowed for the potential discovery of entirely different clusters. That we identified 7 of the 8 similar disease states and all 4 trajectories strongly confirms the robustness and generalisability of the identified structure.

Previous studies on RA trajectories focused on recovery speed and identified factors like age, sex, and baseline ESR as influential [19]. Our findings provide more nuanced insights, revealing elevated ESR vs joint inflammation as characteristics of distinct patient subgroups. We observed both favourable and unfavourable trajectories within seropositive cohorts, and symptom duration varied among trajectories across data sets. These results challenge established notions about predictive factors in RA outcomes and offer a more detailed perspective on disease progression. The divergent disease states at 1.5 years highlight potential areas for clinical improvement, with over 50% of patients ending with elevated ESR or leucocytes. This co-occurs with decreased Hb and elevated thrombocytes within the limits of normal. This suboptimal outcome is mirrored in PROMs and suggests that blood inflammation is not merely a bystander effect.

While the distinction between systemic and localised inflammation emerges as an overarching pattern, it represents a simplification of complex inflammatory processes that likely interact and overlap. This framework, derived from unsupervised clustering rather than imposed a priori hypothesis, provides a clinically interpretable foundation for understanding disease heterogeneity while acknowledging that the underlying biology is more nuanced.

The different patterns observed in the various trajectories may indicate the existence of multiple RA subtypes. Further studies are required to substantiate this hypothesis. The main disease states, trajectories A and B, show a consistent difference in having either predominantly elevated ESR or joint inflammation. Women, older people, and people with aCCP were more likely to end up in trajectory A than those in B. Possibly, the A subtype is more B cell driven, while the B subtype is more driven by innate immunity. This is supported by the higher prevalence

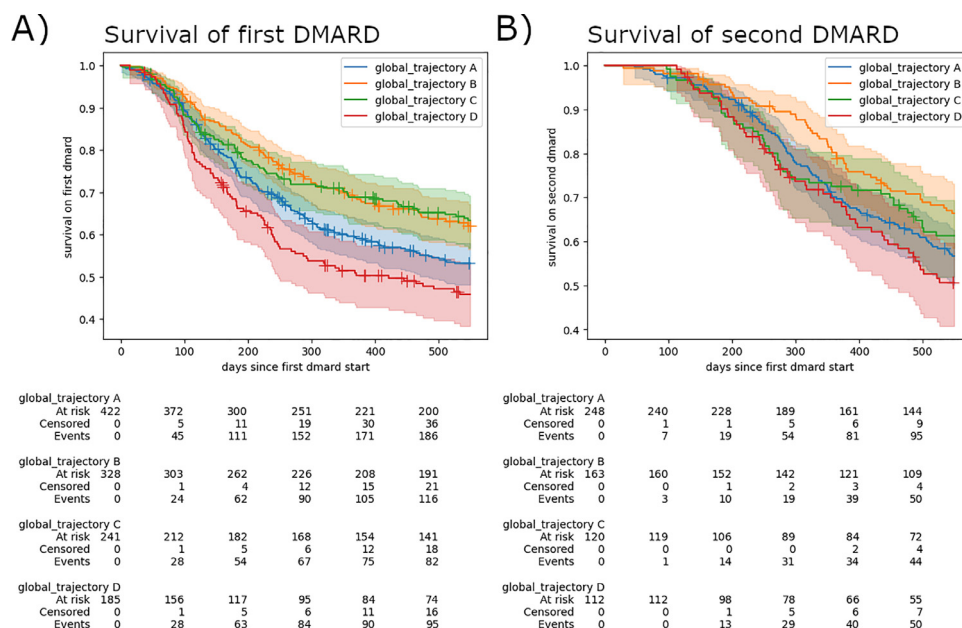


Figure 6. Kaplan–Meier curves for disease-modifying antirheumatic drug (DMARD) survival across disease trajectories. (A) First DMARD survival and (B) second DMARD survival over 500 days from first DMARD initiation. Shaded areas represent 95% CIs. Tables below each plot show the number of patients at risk, censored cases, and events at 100-day intervals. Trajectory B shows superior drug survival for both first and second DMARDs, while trajectory D shows accelerated discontinuation rates.

of seropositivity within the A subtype. Although we have named the inflammation in trajectory A systemic inflammation, it is possible that the elevated ESR is an inflammation from elsewhere in the body, for instance in the lungs or gums, both of which are known to be correlated to aCCP-positive disease [20,21].

The leucocytosis in trajectory C suggests an IL-8–driven subtype. We cannot definitively conclude that trajectory C is not caused by prednisone prescriptions, as we lack data on prednisone use before baseline, although our analyses show high leucocytes in trajectory C patients with and without prednisone.

Trajectory D resembles the pauci-immune or fibroid subtypes described in recent biopsy studies [22,23]. Patients in trajectory D show low ESR, high SJC and TJC early in the disease, low seropositivity, poor treatment response, and a discrepancy between PROMs and physician global assessment. Despite differences in study populations, these features align with suggestions of aCCP-negative RA being more pauci-immune [24,25]. In the seropositive TACERA cohort, trajectory D patients still display these features, suggesting that also within this group, a pauci-immune-like subset might exist. Further study to the 54% patients from trajectory D, who at start resemble the more favourable outcome group of trajectory B, might underpin important treatment opportunities for this difficult-to-treat group.

By excluding previously established predictors (age, sex, and seropositivity) from our clustering algorithm, we could independently validate their relationship with the identified trajectories. Analysing how established predictors like RF status and sex related to our identified trajectories, provided valuable validation of our approach. For example, we found that women and seropositive patients were indeed more likely to follow trajectory A, confirming the relevance of these factors, while demonstrating that they emerge naturally from inflammation-based clustering. Importantly, our analysis also revealed that traditional subset markers may be imperfect proxies for underlying disease patterns, as we observed women and seropositive patients across all trajectories. This suggests that our inflammation-based clustering approach might capture more fundamental disease mechanisms that are only partially reflected by conventional patient characteristics. This observation aligns with the growing recognition that RA is a heterogeneous disease with complex pathophysiological patterns that may not perfectly align with traditional patient subsets. Finally, finding the same clusters within the completely seropositive replication data set of TACERA, confirms our choice of not including serology.

The differential patterns of DMARD survival across trajectories provide additional validation of our clustering approach. The poor DMARD survival in trajectory D aligns with its overall worse outcomes, while the superior persistence in trajectory B supports its identification as a favourable disease course. Notably, the consistency of these patterns across both first and second DMARDs suggests these trajectories represent fundamental disease characteristics rather than simply treatment response categories. The similar DMARD survival between trajectories A and D, despite their different inflammatory patterns (systemic vs local), suggests that both forms of persistent inflammation may similarly impact treatment durability. These findings suggest that early trajectory identification might help inform expectations about treatment durability and could guide decisions about treatment intensity.

Our study has several limitations. First, trajectories may be influenced by deviations from treatment guidelines, leading to indication bias, especially for patients resembling trajectory A

or B but later shifting to C or D. However, 48% and 26% of patients were correctly predicted to follow trajectories C and D at baseline, respectively. While these findings enhance our understanding of disease patterns, the prediction model's clinical utility is limited when using only baseline features. Optimising the algorithm for practical implementation remains a goal for future research beyond this study's scope. Another possible limitation is phenotypical misclassification, such as by relying on rheumatologists' diagnoses in the Leiden data. Despite this, the same trajectories were observed in TACERA, where all patients met RA criteria at baseline.

The third key principle of the feature selection was to exclude medications from the clustering algorithm, with the aim of avoiding the creation of clusters that primarily reflect treatment decisions rather than underlying disease phenotypes. The incorporation of medication data has the potential to introduce bias from unknown factors influencing prescription patterns and treatment alterations over time. The approach adopted in this study separates phenotype identification from treatment variables, enabling subsequent analysis of differential treatment responses across the identified clusters. This, in turn, has the potential to reveal how distinct patient subgroups respond to various interventions. This methodologic decision is one of the cornerstones of the development of more targeted, personalised treatment strategies based on objectively defined disease phenotypes.

There are missing disease states between patients' first visits and the 1.5-year cutoff, which may result in an incomplete or incorrect trajectory. The methodology used partly overcomes this problem. The sequence similarity approach inherently accommodates temporal irregularities by focusing on the pattern of disease states rather than the precise timing of transitions. This results in the possibility that while patients might have a missing visit, the sequence is still similar in the comparison made.

The unsupervised clustering design could present a limitation. More than 4 trajectories may exist, as clustering depends on the number of clusters chosen. In the Leiden data, 8 disease states and 4 trajectories were optimal, and the TACERA analysis was aligned to maintain comparability. Forcing 8 clusters in TACERA resulted in a new cluster, T9, likely a mix of disease states L4 and L7. This suggests that with more extensive data, additional trajectories could emerge, as seen with substantial overlap in trajectory C with trajectories A and B. Trajectory C might even be subdivided further.

In conclusion, our study identified the existence of 4 distinct disease trajectories in early RA, categorised into patterns of localised and systemic inflammation. These trajectories were replicated, demonstrating the robustness of our findings. The results provide a more nuanced understanding of RA progression, challenging traditional predictors and suggesting the presence of multiple RA subtypes. While systemic inflammation was associated with poorer outcomes, localised inflammation also presented challenges, particularly in trajectory D, which may align with the pauci-immune subtype. Future studies should explore further subdivision of trajectories and potential treatment strategies for difficult-to-treat subgroups, particularly those transitioning from more favourable to less favourable outcomes.

Competing interests

All authors declare they have no competing interests.

CRedit authorship contribution statement

Nils Steinz: Writing – review & editing, Writing – original draft, Visualization, Validation, Software, Methodology, Investigation, Formal analysis, Data curation. **Tjardo D. Maarseveen:** Writing – review & editing, Validation, Methodology, Data curation. **Erik B. van den Akker:** Writing – review & editing, Visualization, Validation, Supervision, Methodology. **Andrew P. Cope:** Writing – review & editing, Funding acquisition, Conceptualization. **John D. Isaacs:** Writing – review & editing, Funding acquisition, Data curation, Conceptualization. **Aaron R. Winkler:** Writing – review & editing, Funding acquisition, Conceptualization. **Tom W. J. Huizinga:** Writing – review & editing, Supervision, Resources, Project administration, Funding acquisition. **Yann Abraham:** Writing – review & editing, Supervision, Methodology, Investigation, Funding acquisition, Data curation, Conceptualization. **Rachel Knevel:** Writing – review & editing, Writing – original draft, Visualization, Validation, Supervision, Project administration, Methodology, Investigation, Funding acquisition, Formal analysis, Data curation, Conceptualization.

Acknowledgements

We thank Samantha Jurado-Zapata and David Steeman for their help in extracting and processing the electronic health record data from Leiden University Medical Centre. Additional thanks to Daan van der Bijl for making the first steps in the study. The results in this paper were presented at EULAR 2024 under abstract number OP0018 (doi:10.1136/annrheumdis-2024-eular.1512).

Funding

This project has received funding from Horizon Europe programme [101095052] (SQUEEZE), [101080711] (SPIDERR), and [777357] (RTCure); MRC/ABPI Inflammation and Immunology Initiative Grant [MRC reference numbers: G1001516 and G1001518] Towards a Cure for Early Rheumatoid Arthritis (TACERA); ZonMw klinische fellow [40-00703-97-19069]; and ZonMw Open Competitie 2021 [09120012110075].

Patient consent for publication

Not applicable.

Ethics approval

Ethics committee approved.

Provenance and peer review

Not commissioned; externally peer reviewed

Data availability statement

We have made our scripts available in a public repository at: https://github.com/nilssteinz/Early_RA_Trajectories. Study data are available upon reasonable request.

Declaration of generative AI and AI-assisted technologies in the writing process

During the preparation of this work, the authors used Claude.ai in order to improve the text and the readability of the text. After using this tool/service, the authors reviewed and edited the content as needed and take(s) full responsibility for the content of the published article.

Supplementary materials

Supplementary material associated with this article can be found in the online version at doi:10.1016/j.ard.2025.04.011.

Orcid

Nils Steinz: <http://orcid.org/0009-0009-4314-8228>

Erik B. van den Akker: <http://orcid.org/0000-0002-7693-0728>

John D. Isaacs: <http://orcid.org/0000-0002-6103-7056>

Aaron R. Winkler: <http://orcid.org/0000-0001-7365-795X>

Yann Abraham: <http://orcid.org/0000-0001-5600-5896>

REFERENCES

- [1] Watanabe R, Hashimoto M, Murata K, Murakami K, Tanaka M, Ohmura K, et al. Prevalence and predictive factors of difficult-to-treat rheumatoid arthritis: the KURAMA cohort. *Immunol Med* 2022 Jan;45(1):35–44.
- [2] Ni J, Wang P, Yin KJ, Yang XK, Cen H, Sui C, et al. Novel insight into the aetiology of rheumatoid arthritis gained by a cross-tissue transcriptome-wide association study. *RMD Open* 2022 Sep;8(2):e002529.
- [3] van der Kooij SM, Goekoop-Ruiterman YPM, de Vries-Bouwstra JK, Peeters AJ, van Krugten MV, Breedveld FC, et al. Probability of continued low disease activity in patients with recent onset rheumatoid arthritis treated according to the disease activity score. *Ann Rheum Dis* 2008 Feb;67(2):266–9.
- [4] Studenic P, Radner H, Smolen JS, Aletaha D. Discrepancies between patients and physicians in their perceptions of rheumatoid arthritis disease activity. *Arthritis Rheum* 2012;64(9):2814–23.
- [5] Singh JA, Saag KG, Bridges Jr SL, Akl EA, Bannuru RR, Sullivan MC, et al. 2015 American College of Rheumatology guideline for the treatment of rheumatoid arthritis. *Arthritis Rheumatol* 2016;68(1):1–26.
- [6] Movahedi M, Cesta A, Li X, Bombardier C, investigators OBRI. Disease activity trajectories for early and established rheumatoid arthritis: real-world data from a rheumatoid arthritis cohort. *PLoS One* 2022 Sep;17(9):e0274264.
- [7] RA-MAP Consortium. Characterization of disease course and remission in early seropositive rheumatoid arthritis: results from the TACERA longitudinal cohort study. *Ther Adv Musculoskelet Dis* 2021; 13:1759720X211043977.
- [8] Barnabe C, Sun Y, Boire G, Hitchon CA, Haraoui B, Thorne JC, et al. Heterogeneous disease trajectories explain variable radiographic, function and quality of life outcomes in the Canadian Early Arthritis Cohort (CATCH). *PLoS One* 2015 Aug;10(8):e0135327.
- [9] Ahlqvist E, Storm P, Käräjämäki A, Martinell M, Dorkhan M, Carlsson A, et al. Novel subgroups of adult-onset diabetes and their association with outcomes: a data-driven cluster analysis of six variables. *Lancet Diabetes Endocrinol* 2018 May;6(5):361–9.
- [10] Ma EY, Kim JW, Lee Y, Cho SW, Kim H, Kim JK. Combined unsupervised-supervised machine learning for phenotyping complex diseases with its application to obstructive sleep apnea. *Sci Rep* 2021 Feb;11(1):4457.
- [11] Maurits MP, Korsunsky I, Raychaudhuri S, Murphy SN, Smoller JW, Weiss ST, et al. A framework for employing longitudinally collected multicenter electronic health records to stratify heterogeneous patient populations on disease history. *J Am Med Inform Assoc* 2022 May;29(5):761–9.
- [12] Heumos L, Schaar AC, Lance C, Litinetskaya A, Drost F, Zappia L, et al. Best practices for single-cell analysis across modalities. *Nat Rev Genet* 2023 Aug;24(8):550–72.
- [13] Maarseveen TD, Maurits MP, Niemantsverdriet E, van der Helm-van Mil AHM, Huizinga TWJ, Knevel R. Handwork vs machine: a comparison of rheumatoid arthritis patient populations as identified from EHR free-text by diagnosis extraction through machine-learning or traditional criteria-based chart review. *Arthritis Res Ther* 2021 Jun;23(1):174.

- [14] Arnett FC, Edworthy SM, Bloch DA, Mcshane DJ, Fries JF, Cooper NS, et al. The American Rheumatism Association 1987 revised criteria for the classification of rheumatoid arthritis. *Arthritis Rheum* 1988;31(3):315–24.
- [15] Aletaha D, Neogi T, Silman AJ, Funovits J, Felson DT, Bingham III CO, et al. 2010 Rheumatoid arthritis classification criteria: an American College of Rheumatology/European League Against Rheumatism collaborative initiative. *Arthritis Rheum* 2010 Sep;62(9):2569–81.
- [16] Deighton C, O'Mahony R, Tosh J, Turner C, Rudolf M, Guideline Development Group. Management of rheumatoid arthritis: summary of NICE guidance. *BMJ* 2009 Mar;338:b702.
- [17] Levine JH, Simonds EF, Bendall SC, Davis KL, el-AD Amir, Tadmor MD, et al. Data-driven phenotypic dissection of AML reveals progenitor-like cells that correlate with prognosis. *Cell* 2015 Jul;162(1):184–97.
- [18] Bobroske K, Larish C, Cattrell A, Bjarnadóttir MV, Huan L. The bird's-eye view: a data-driven approach to understanding patient journeys from claims data. *J Am Med Inform Assoc* 2020 Jun;27(7):1037–45.
- [19] Nilsson J, Andersson MLE, Hafström I, Svensson B, Forslind K, Ajeganova S, et al. Influence of age and sex on disease course and treatment in rheumatoid arthritis. *Open Access Rheumatol* 2021;13:123–38.
- [20] Holers VM, Demoruelle MK, Kuhn KA, Buckner JH, Robinson WH, Okamoto Y, et al. Rheumatoid arthritis and the mucosal origins hypothesis: protection turns to destruction. *Nat Rev Rheumatol* 2018 Sep;14(9):542–57.
- [21] Demoruelle MK, Deane KD, Holers VM. When and where does inflammation begin in rheumatoid arthritis? *Curr Opin Rheumatol* 2014 Jan;26(1):64–71.
- [22] Zhang F, Jonsson AH, Nathan A, Millard N, Curtis M, Xiao Q, et al. Deconstruction of rheumatoid arthritis synovium defines inflammatory subtypes. *Nature* 2023 Nov;623(7987):616–24.
- [23] Matthijssen XME, Niemantsverdriet E, Huizinga TWJ, van der Helm–van Mil AHM. Enhanced treatment strategies and distinct disease outcomes among autoantibody-positive and -negative rheumatoid arthritis patients over 25 years: a longitudinal cohort study in the Netherlands. *PLoS Med* 2020 Sep;17(9):e1003296.
- [24] Nerviani A, Di Cicco M, Mahto A, Lliso-Ribera G, Rivellesse F, Thorborn G, et al. A pauci-immune synovial pathotype predicts inadequate response to TNF α blockade in rheumatoid arthritis patients. *Front Immunol* 2020;11:845.
- [25] Li K, Wang M, Zhao L, Liu Y, Zhang X. ACPA-negative rheumatoid arthritis: from immune mechanisms to clinical translation. *EBioMedicine* 2022 Sep;83:104233.