

# Forecasting Clinical Outcomes from Limited Real-World Data: A Comparative Simulation Study

Awa Diop<sup>1</sup>, Sheena Kayaniyil<sup>1</sup>, Lise Retat<sup>2</sup>, Sarah Collier<sup>1</sup>, Diar Fattah<sup>3</sup>, Stefan Franzén<sup>4,5</sup>

<sup>1</sup>BioPharmaceuticals Medical, AstraZeneca, Mississauga, Canada | <sup>2</sup>Global Market Access and Pricing, AstraZeneca, Barcelona, Spain | <sup>3</sup>BioPharmaceuticals Medical, AstraZeneca, Barcelona, Spain  
<sup>4</sup>BioPharmaceuticals Medical, AstraZeneca, Gothenburg, Sweden | <sup>5</sup>School of Public Health and Community Medicine, Institute of Medicine, Sahlgrenska Academy, University of Gothenburg, Sweden

## BACKGROUND

- AstraZeneca's location-relevant real-world data (RWD) platform, ATLAS, captures retrospective longitudinal healthcare data to enable country-level informed decision across multiple disease areas.
- Estimating future trends in healthcare outcomes at national and sub-national levels is essential for effective local planning and can uncover opportunities where interventions can influence meaningful improvements.
- Using RWD to project clinical outcomes can be effective because it can be accessible, timely, robust, and representative across small geographies.
- Traditional forecasting methods struggle with structural breaks (e.g., COVID-19), compounded by limited historical location-specific RWD and disease-specific seasonality impacts.

## OBJECTIVES

- Compare modeling approaches (indicated below) for forecasting key clinical outcomes from limited RWD.
- Assess robustness and short-horizon performance of:
  - Linear Model (LM)
  - Generalized Additive Model (GAM)
  - Autoregressive Integrated Moving Average (ARIMA), each combined with Interrupted Time Series (ITS) to capture COVID-related disruptions.<sup>1,2</sup>
- Evaluate separate vs pooled vs hierarchical analyses of sub-national regions.

## METHODS

- Study design:** Simulation study.
- Synthetic RWD data was generated to inform a simulation process. COVID-19 disruption was considered as the middle of available time points (as depicted in **Figure 1**).
- Available RWD:**
  - Chronic obstructive pulmonary disease (COPD) hospital admission rates & counts.
  - One national series and 18 regional series considering both, with and without seasonality (using Spain real-world healthcare data as a framework)
  - T = 10 to 50 time points of follow-up (in synthetic RWD)
- Outcomes:** hospital admission rates and counts
- Evaluation:** Rolling-origin cross-validation at future time points: 1, 3, 5 (horizon)
- Time points are based on data availability (e.g. can be monthly, quarterly, annually)
- Analyses were conducted to investigate the accuracy and uncertain of the models, with and without seasonality considerations
- Metrics:**
  - Accuracy: RMSE (Root Mean Squared Error), MAE (Mean Absolute Error), sMAPE (Symmetric Mean Absolute Percentage Error)
  - Uncertainty: 95% PI Coverage, PI Width

Key design feature: ITS component explicitly models the COVID-19 structural break, allowing all three base models to account for pandemic disruptions.

## DATA-GENERATING PROCESS

Figure 1. Synthetic RWD presented as (A) Non-seasonal time series (B) Seasonal time series

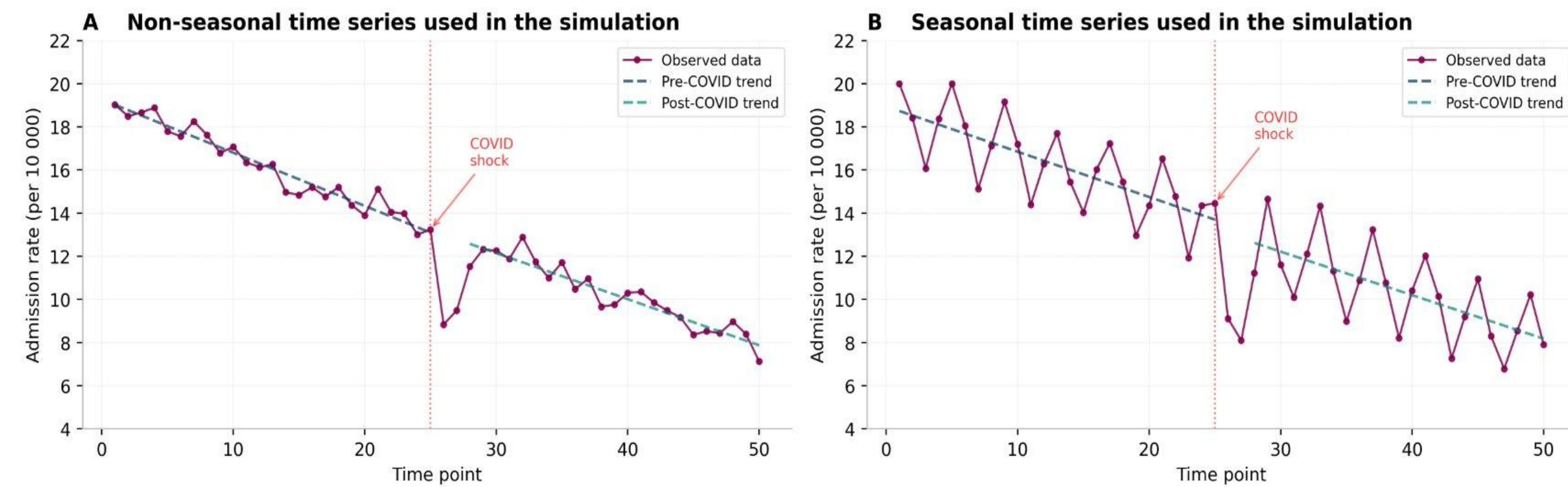


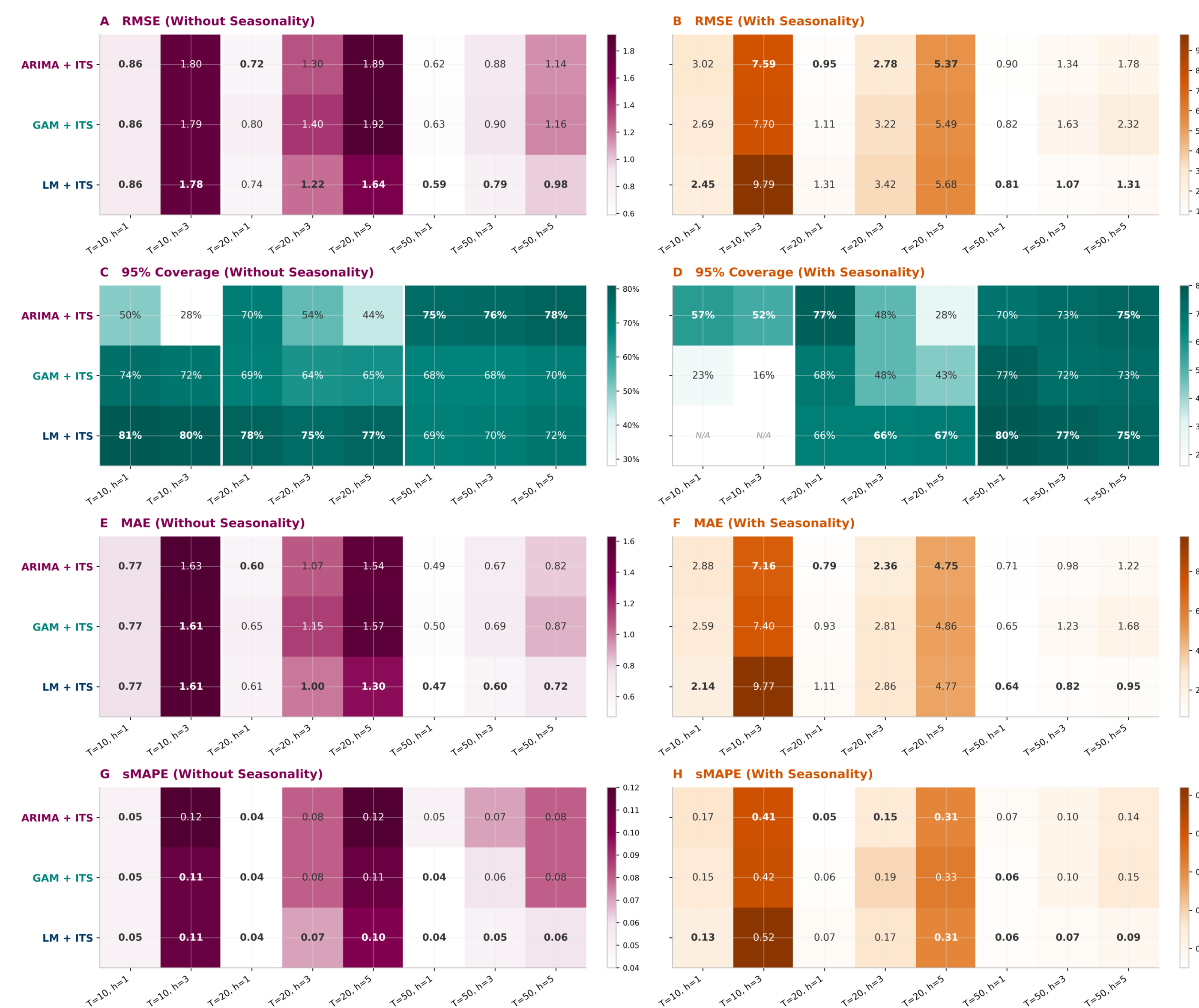
Figure 1A and 1B describe the data-generating process for this simulation study.

The simulation is based on RWD observed for admissions rate and count to avoid arbitrary parameter choices and allow fair method comparisons.

## RESULTS

40 scenarios were evaluated, as depicted in **Figure 2**.

Figure 2. National level results: (A) RMSE without seasonality; (B) RMSE with seasonality; (C) 95% coverage without seasonality; (D) 95% coverage with seasonality; (E) MAE without seasonality; (F) MAE with seasonality; (G) sMAPE without seasonality; (H) sMAPE with seasonality

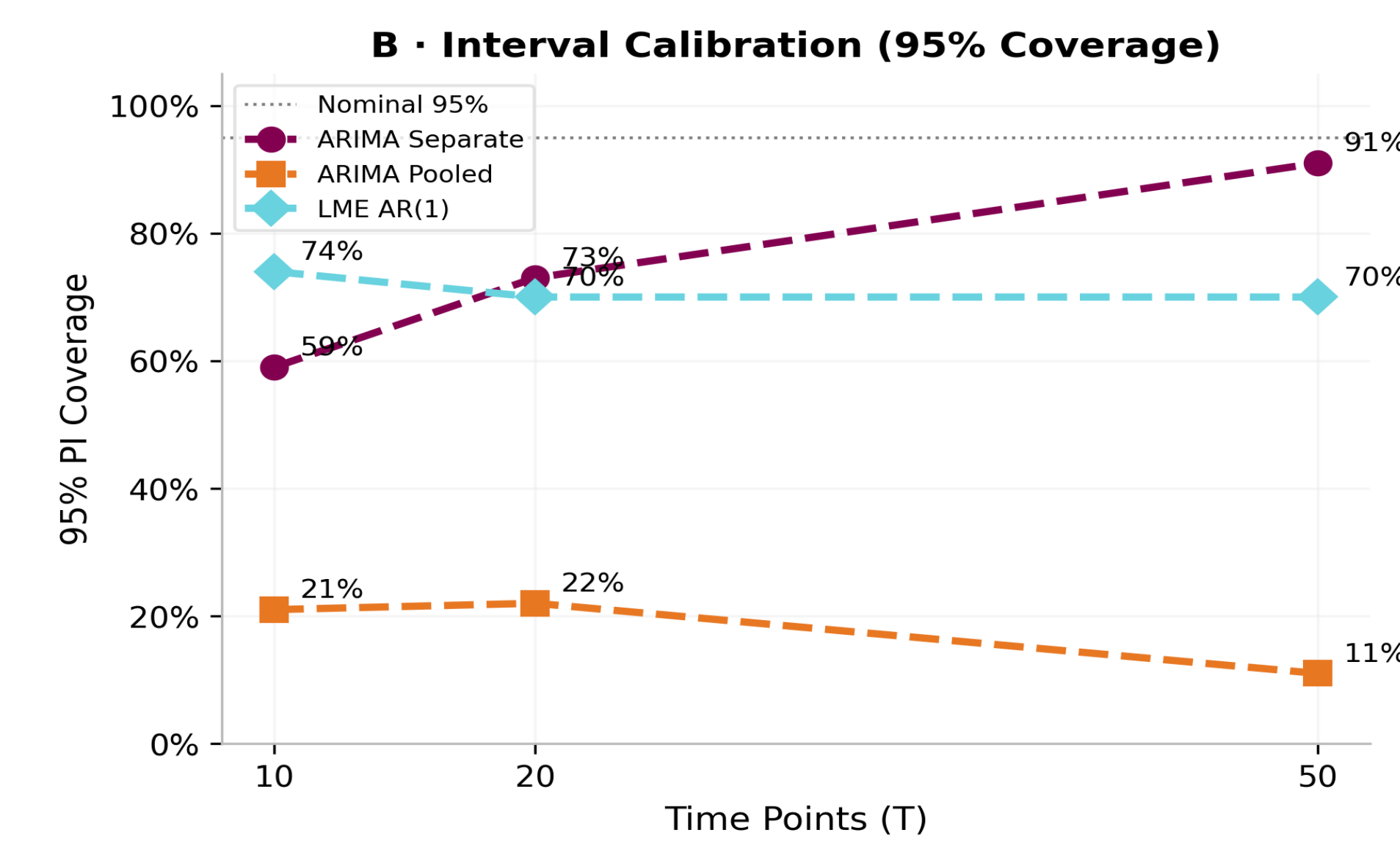
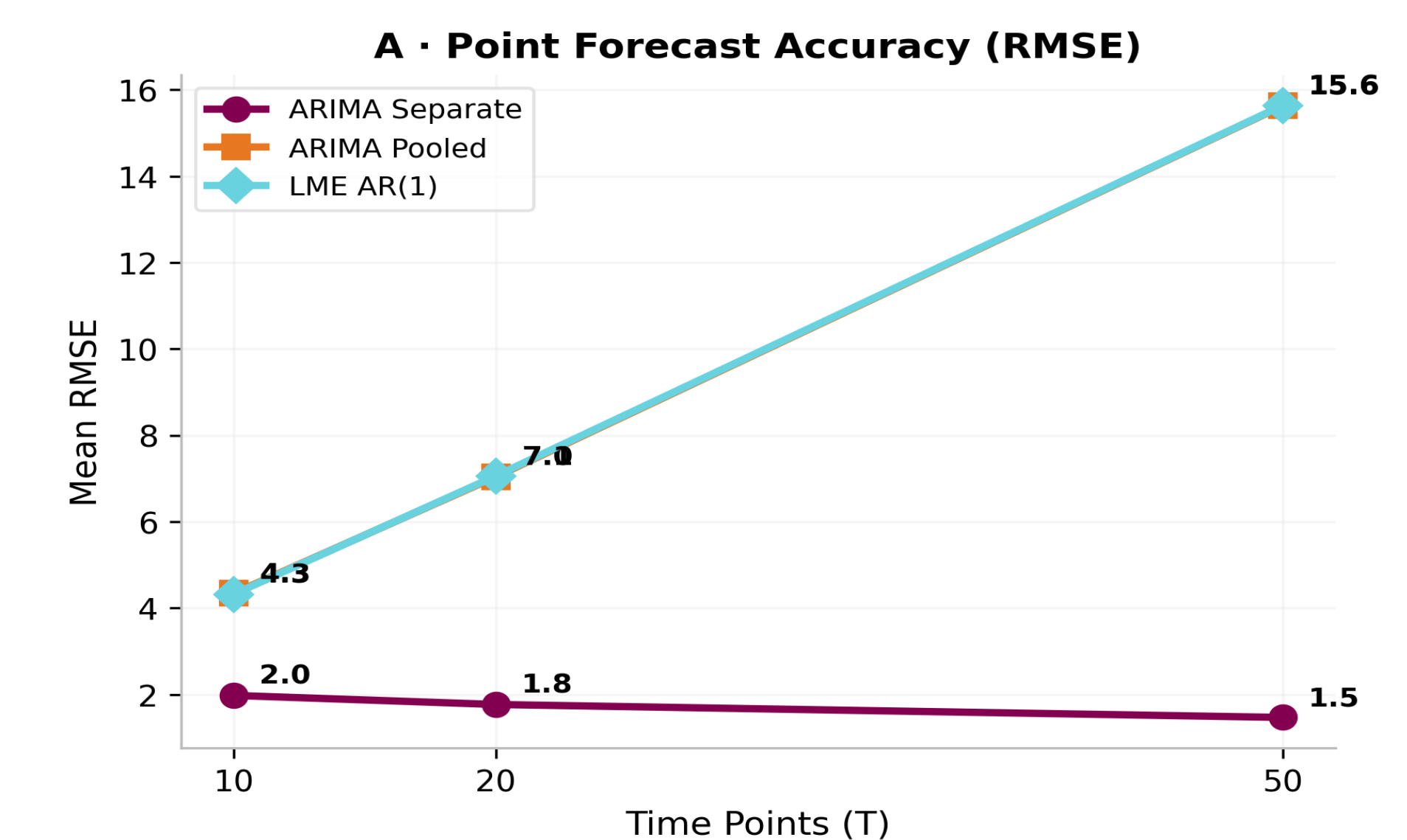


How to read: Each row shows a metric; left column = without seasonality, right column = with seasonality. Each cell shows a model (row) vs a given scenario (column). For time points, it is forecast horizon. Bold values = best model per scenario. N/A = model did not converge.

Without seasonality, across 10-50 time points, LM + ITS, GAM + ITS, and ARIMA + ITS achieved similar accuracy (e.g., for T = 10 and 3 future time points, RMSE were respectively 1.78, 1.79 and 1.80), **Figure 2A**.

With seasonality, ARIMA+ITS achieved overall lower accuracy (RMSE), **Figure 2B**.

Figure 3. Regional analysis – ARIMA + ITS applied separately in each region (“ARIMA Separate”), ARIMA + ITS applied to all regions combined (“ARIMA Pooled”) and Hierarchical (“LME + AR (1)”: (A) RMSE; (B) 95% coverage.



LME + AR (1): Linear Mixed-effects (LME) model where the residual errors follow an AR(1) autocorrelation structure

Figure 3 compares the regional analysis as a separated, pooled or hierarchical approach (considers correlation between regions).

The ARIMA + ITS Separate model outperform pooled and hierarchical approaches in RMSE **Figure 3A**.

Accounting for correlations between regions improved uncertainty (95% coverage) more than improving RMSE **Figure 3B**.

## KEY INSIGHTS

- Without seasonality: all three ITS models are interchangeable for point prediction.
- With seasonality: ARIMA + ITS demonstrated superior performance, especially with shorter time points in data.
- Regional forecasting: ARIMA + ITS applied for reach region separately delivers the best point accuracy across all time points of observed data.
- Pooling across regions helps uncertainty, not point error — hierarchical fits improve interval calibration when T is large, but adds complexity without reducing RMSE.

## CONCLUSIONS

- ARIMA + ITS was found to be a flexible and reliable approach for timely, informed decision-making. It adapts to different data configurations and can account for seasonality and structural interruptions in the time series.
- Separate region-by-region fitting is preferred for heterogeneous panels.
- Even with shorter time horizons, pooling/hierarchical structures is preferable for uncertainty, yet performed poorly for accuracy, compared to ARIMA Separate.
- ITS was successful in correctly handling COVID-era disruptions in health outcome data.
- A key limitation is the synthetic RWD based on one disease and one healthcare metric. Yet, we incorporated various follow-up time points and potential seasonality impacts.

## REAL WORLD APPLICATION

This work will directly inform AstraZeneca's location-relevant RWD platform, ATLAS, and will be considered for additional markets and outcomes.<sup>3</sup> Generating forecasts for key outcomes of interest is critical for planning purposes and can enable an urgency to act now to minimize future burden. Demonstrating region-specific future trends, which have accounted for COVID-19 disruptions, is valuable to unlock value for healthcare decision making.

## REFERENCES

- Schaffer AL, et al. Interrupted time series analysis using ARIMA models: a guide for evaluating large-scale health interventions. BMC Med Res Methodol. 2021;21:58.
- Zhou, Q., Hu, J., Hu, W., Li, H., & Lin, G. Z. (2023). Interrupted time series analysis using the ARIMA model of the impact of COVID-19 on the incidence rate of notifiable communicable diseases in China. BMC Infectious Diseases, 23(1), 375.
- Medine J, Sánchez-Covisa, J, Nuevo J, et al. (2026) Shaping the future of evidence generation: Real-world data to drive healthcare transformation and patient-centered decisions.

Conflict of interest: AD, SK, LR, SC, DF and SF are all employees of AstraZeneca. AD, SK, SC, DF and SF each have stock ownership and/or stock options or interests in AstraZeneca.