

# Millennial Medical Record Data Profile: A Japanese Electronic Medical Records Database Utilising Unstructured Data for Lung Cancer Research

Ayumi Hamaguchi,<sup>1</sup> Kenichiro Shiba,<sup>2</sup> Mari Kato,<sup>3</sup> Daisuke Wakasugi,<sup>4</sup> Tadashi Koga,<sup>3</sup> Suguru Nozue,<sup>4</sup> Takayuki Sawada,<sup>3</sup> Tomoyo Morita,<sup>4</sup> Amanda Pulfer,<sup>1</sup> Dimitra Lambrelli<sup>1</sup>

<sup>1</sup>Thermo Fisher Scientific, London, UK; <sup>2</sup>General Incorporated Association Life Data Initiative, Kyoto, Japan; <sup>3</sup>Clinical Study Support, Inc., Nagoya, Japan; <sup>4</sup>NTT DATA Japan Corporation, Tokyo, Japan

## Background

- Real-world data sources, such as health insurance claims and Diagnosis Procedure Combination (DPC) databases, are commonly used in Japanese research because of their broad coverage and suitability for longitudinal analyses,<sup>1,2</sup> but these sources often lack detailed clinical information.
- Electronic medical records (EMRs) offer rich clinical detail, including laboratory values and indicators of disease severity and progression,<sup>3</sup> but access to EMR-based data—particularly for commercial research—remains limited.
- The Millennial Medical Record (MMR) database includes data from EMRs, claims, and DPC sources. It has the potential to be leveraged for real-world research, including lung cancer and other disease areas, using its comprehensive clinical and longitudinal data.
- To describe the characteristics of the MMR database, lung cancer was selected due to its clinical relevance and its status in Japan as a leading cause of death,<sup>4</sup> with its rising age-adjusted incidence rate since 2019.<sup>5</sup>

## Objectives

- The goal of this study was to describe the characteristics of MMR data, illustrating its utility for lung cancer research.

## Methods

### MMR Overview

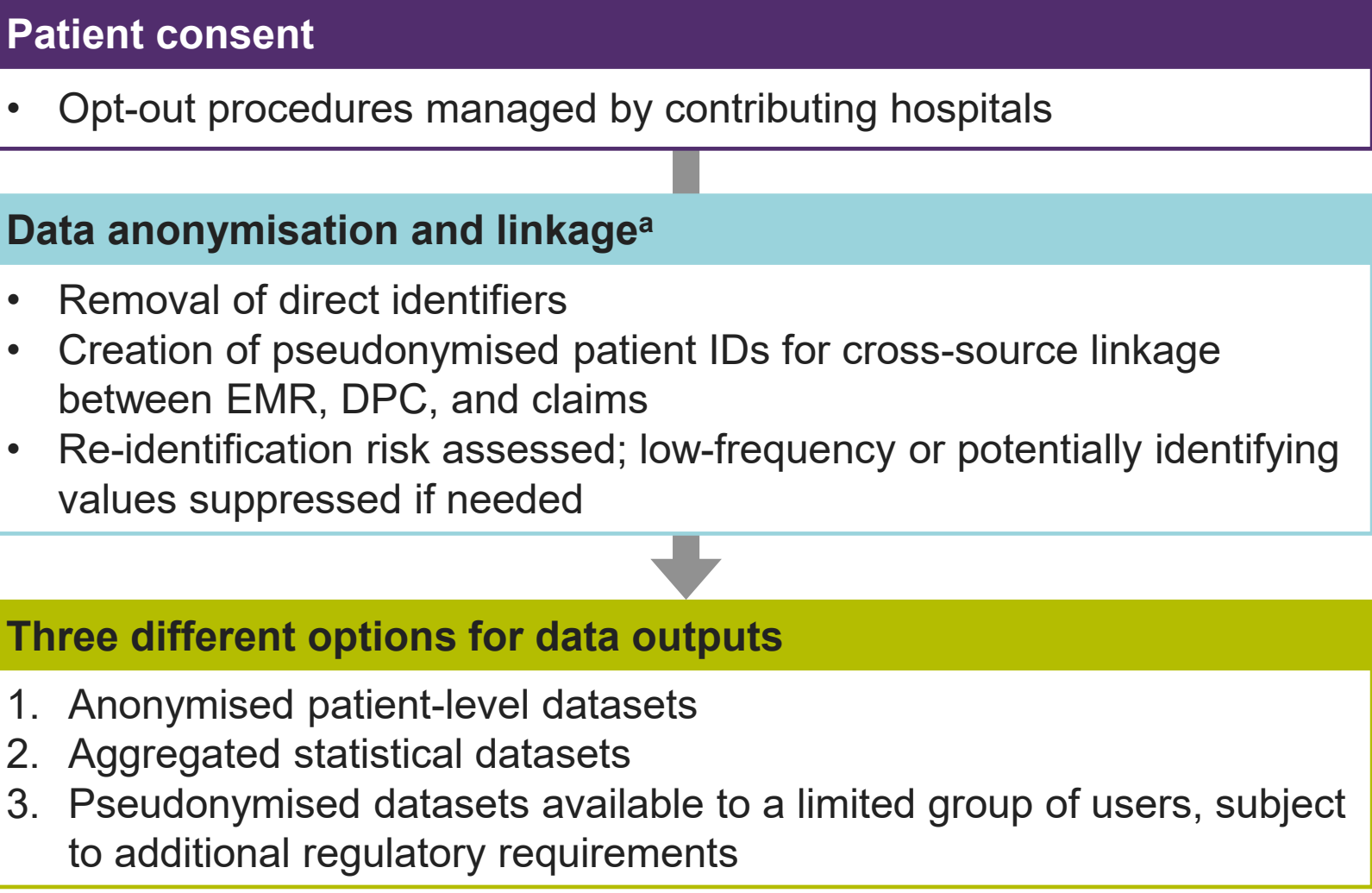
- The MMR database includes structured and unstructured data from EMRs (**Table 1**) that complies with the Medical Markup Language 4.2.0 standard,<sup>6</sup> claims data, DPC data, and imaging data.
- Data are collected under the framework of the Next Generation Medical Infrastructure Act,<sup>7</sup> which allows only certified entities to manage and anonymise medical information for research use. Processing staff complete mandatory training in anonymisation, data security, and compliance (**Figure 1**). The MMR database is operated by the Life Data Initiative and NTT DATA Japan Corporation.
- As MMR's coverage expands throughout Japan, the number of participating hospitals increased from 43 in 2020 to 67 in 2025 (**Figure 2**). Of these, the number contributing EMR, DPC, and claims data usable for research increased from 10 to 26 over the same period.
- Figure 3** provides an overview of data availability in the MMR database.

Table 1. MMR Overview

Participating hospitals	67 secondary/tertiary care hospitals with an average of 600 beds, mostly university and regional core hospitals under secondary use agreements with LDI, as of Aug 2025
Patient coverage	1.9 million patients (approx. 1.5% of Japanese population) have complete EMR + DPC + claims data, as of Feb 2025
Data collection	Started in 2016 and is ongoing <sup>a</sup>
Update frequency	Quarterly updates <sup>b</sup>
Data access	Review required by the LDI 'Purpose of Data Use Review Board'

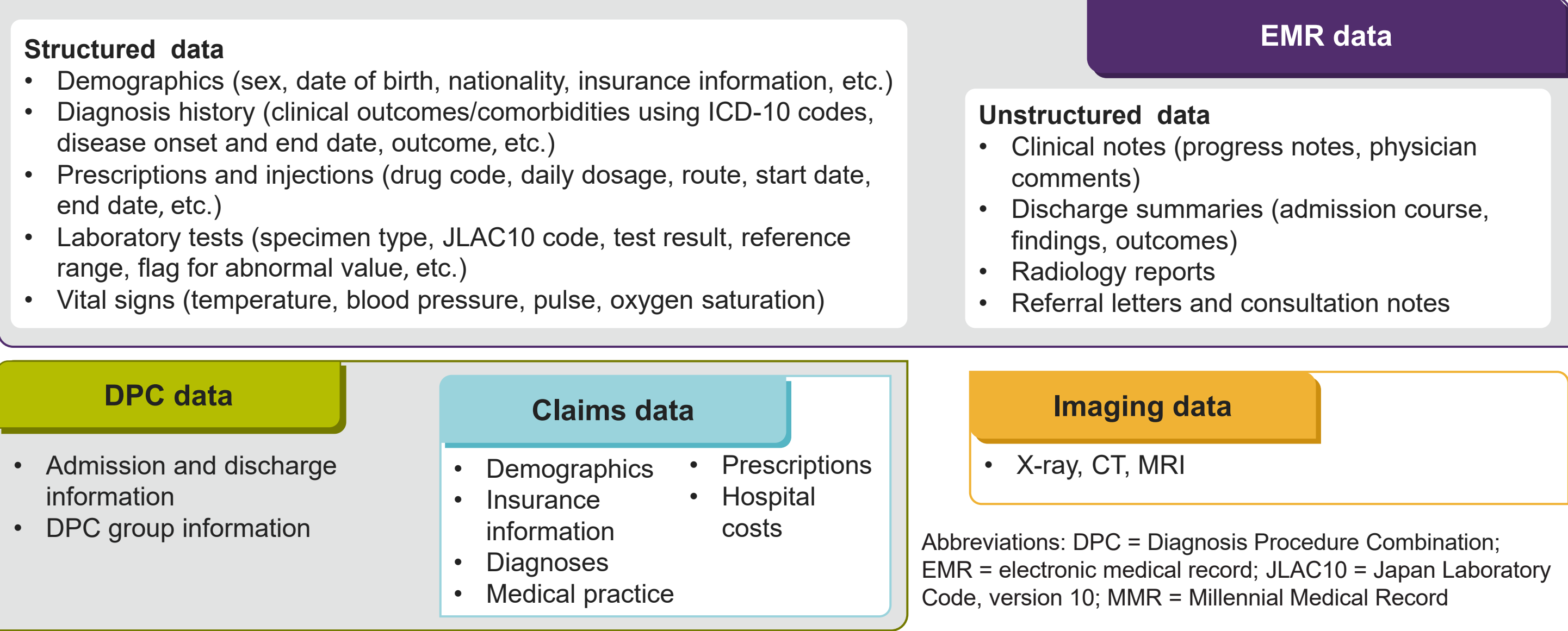
Abbreviations: DPC = Diagnosis Procedure Combination; EMR = electronic medical record; LDI = Life Data Initiative; MMR = Millennial Medical Record.  
<sup>a</sup>The latest data available as of September 2025 is from March 2025. <sup>b</sup>Data updated via automatic transfer or secure file upload from participating hospitals.

Figure 1. Data Processing and Anonymisation



Abbreviations: DPC = Diagnosis Procedure Combination; EMR = electronic medical record; ID = identification.  
<sup>a</sup>Data anonymisation must be performed in a secure workspace (e.g., fingerprint access, camera monitoring, no personal devices).

Figure 3. MMR Data Overview



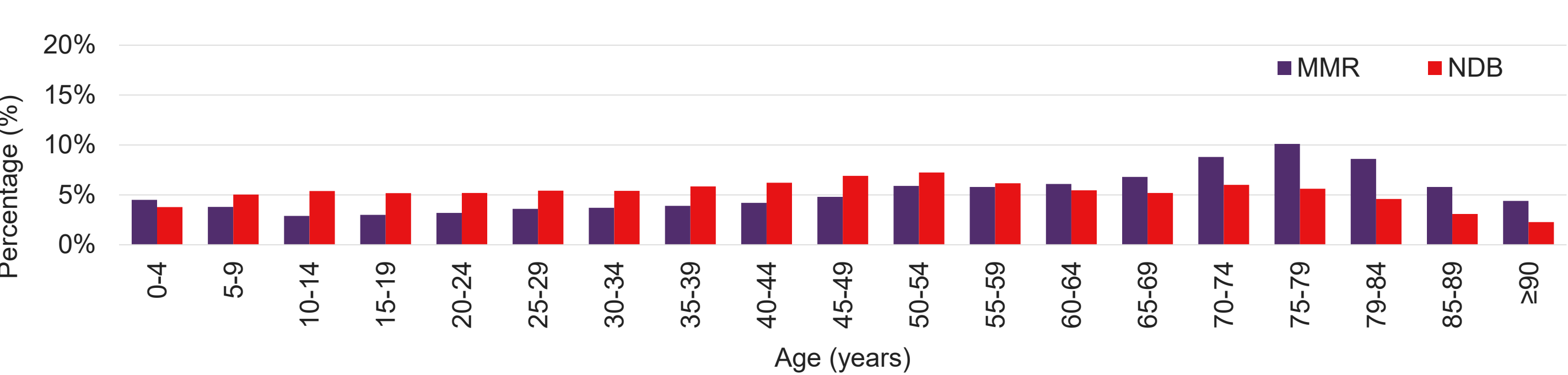
### Statistical Analysis

- To describe MMR data characteristics, age distribution as of 2024 was calculated and compared with the National Database of Health Insurance Claims and Specific Health Checkups of Japan (NDB),<sup>8</sup> using patients with a first-visit claim between April 2023 and March 2024.
- Patients with lung cancer were identified using ICD-10 code a C34\* as of February 2025, with counts stratified by sex and age.
- To demonstrate the value of EMR data, keywords were grouped into five domains: histology, genetic tests, disease scores, metastasis sites, and tumour marker tests. A rule-based keyword search was applied to assess data availability, regardless of confirmed presence. Keywords for histology, genetic tests, and disease scores were extracted from unstructured data using SQL scripts program supported by Generative AI. Analyses were descriptive.

## Results

- As of 2024, the MMR database included 1.6 million individuals from 25 hospitals contributing EMR, DPC, and claims data across Japan, while the NDB database included 88 million individuals. Overall age distributions were similar, with all proportional differences under 5%, supporting the MMR database's representativeness. The largest gap was 4.49% in the 75–79 age group (**Figure 4**).

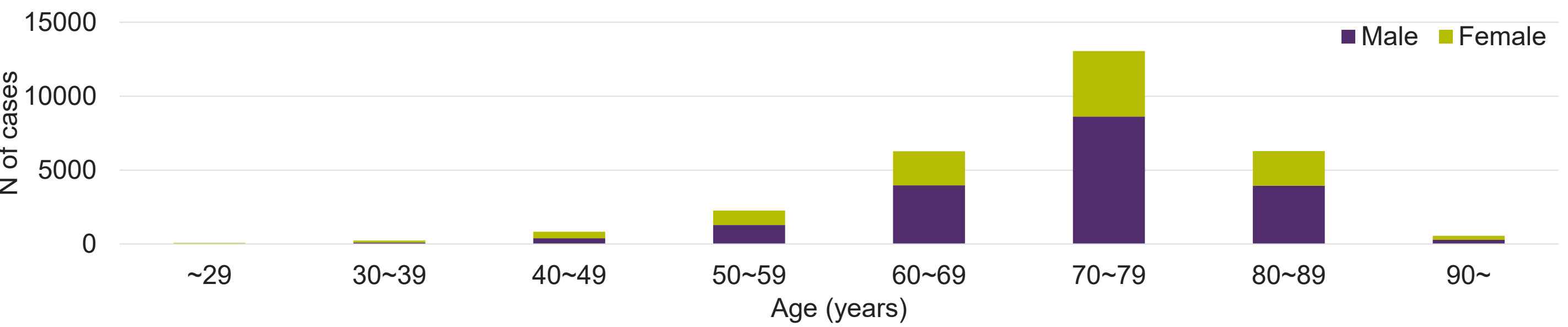
Figure 4. Age Distribution Among MMR and NDB



Abbreviations: MMR = Millennial Medical Record; NDB = National Database of Health Insurance Claims and Specific Health Checkups of Japan.  
Note: The NDB was used as a reference due to its comprehensive nationwide coverage.

- Among 29,545 patients with lung cancer, 63.1% were male (**Figure 5**). The largest age group was 70–79 years old (44.2%), consistent with national incidence data showing higher occurrence in adults aged 69 and older than in younger age groups.<sup>9</sup>

Figure 5. Age and Sex Distribution of Patients With Lung Cancer in MMR



Abbreviation: MMR = Millennial Medical Record

- Among 29,545 patients with lung cancer, keywords indicating tissue types—adenocarcinoma, small cell, and squamous cell carcinoma—were identified in 10%–40% of cases from unstructured data (**Table 2**). The keyword list was indicative; clinician validation would expand data identification.

Table 2. Availability of Clinical Data in Unstructured and Structured EMR

Variables	Keywords Lists/Description	Percentage (%) of Patients With Keywords/Codes
Keywords in unstructured EMR		
Histology		
Adenocarcinoma	(Adenocarcinoma or Ade)+(肺 or Lung)+without negative expressions <sup>a</sup> ; OR (肺)+(腺癌 or 腺がん)+without negative expressions <sup>a</sup>	35~40%
Small cell carcinoma	(Small Cell Lung Cancer or SCLC)+without negative expressions <sup>a</sup> ; OR (肺)+(「非」小細胞)+without negative expressions <sup>a</sup> ; OR (「非」小細胞)+(肺)+without negative expressions <sup>a</sup>	10~15%
Squamous cell carcinoma	(Squamous Cell Carcinoma or SCC)+(肺 or Lung); OR (SCCLC)+without negative expressions <sup>a</sup> ; OR (肺)+(扁平上皮)+without negative expressions <sup>a</sup> ; OR (扁平上皮)+(肺)+without negative expressions <sup>a</sup>	10~15%
Genetic test		
EGFR	(EGFR) <sup>b</sup>	20~25%
ALK	(ALK) <sup>b</sup>	15~20%
ROS1	(ROS1) <sup>b</sup>	5~10%
Disease score		
Performance status	(PS)+Any of the numbers 0–9; OR (Performance)+(Status)+any of the numbers 0–9	55~60%
Metastasis site		
Bone	(骨 or Bone)+(転移 or 病変 or Meta)+without negative expressions <sup>a</sup>	Keywords shown for reference purposes
Lymph node	(リンパ)+(転移 or 病変 or Meta)+without negative expressions <sup>a</sup>	
Brain	(脳 or Brain)+(転移 or 病変 or Meta)+without negative expressions <sup>a</sup>	
Data from structured EMR		
Tumour marker tests		
CEA	–	70~75%
SCC	–	35~40%
CYFRA21-1	–	15~20%

Abbreviations: ALK = anaplastic lymphoma kinase; CEA = carcinoembryonic antigen; CYFRA21-1 = cytokeratin 19 fragment; EGFR = epidermal growth factor receptor; PS = performance status; SCC = squamous cell carcinoma; SCCLC = squamous cell carcinoma lung cancer; SCLC = small-cell lung cancer. <sup>a</sup>Negative expressions indicating absence or uncertainty of a specific condition in Japanese include なし, ない, 無, 疑, ません, or 難し.  
<sup>b</sup>Any alphabet before EGFR/ALK/ROS1 are excluded.

## Conclusions

- Age distributions were generally similar between the MMR and NDB databases, supporting its representativeness.
- By integrating complete EMR data—including histological type and biomarkers from unstructured records—with claims and DPC data, the MMR database enables research closely aligned with routine clinical practice, creating opportunities for real-world studies that complement clinical trials.
- The growth of the MMR database demonstrates the feasibility of harnessing EMR data for large-scale research in Japan. As coverage expands, the MMR will support innovative real-world evidence generation—including external control arms and longitudinal studies—ultimately accelerating advances in treatment and health policy.

### References

1. Sato S, Yasunaga H. *Ann Clin Epidemiol*. 2023;5(2):58-64. 2. Fujinaga J, Fukuoka T. *Drugs - Real World Outcomes*. 2022;9(4):543-550. 3. Zhao Y, Tsubota T. *JMIR Med Inform*. 2023;11:e39876. 4. World Health Organization. Global Cancer Observatory - Japan. 2022. <https://gco.iarc.who.int/media/globocan/factsheets/populations/392-japan-fact-sheet.pdf>. 5. Yamamoto H, et al. *Int J Clin Oncol*. 2025;30(2):199-209. 7. Kobayashi S, Kume N. *J Med Syst*. 2020;44(4):69. 7. Ministry of Health, Labour, and Welfare. Next-generation medical infrastructure law. 2023. <https://www.mhlw.go.jp/content/10808000/001166476.pdf>. In Japanese. 8. Ministry of Health, Labour, and Welfare. The 10th NDB Open Data. [https://www.mhlw.go.jp/stf/seisakunitsuite/bunya/0000177221\\_00016.html](https://www.mhlw.go.jp/stf/seisakunitsuite/bunya/0000177221_00016.html). 9. Foundation for Promotion of Cancer Research. Cancer Statistics in Japan. 2025. [https://ganjoho.jp/public/qa\\_links/report/statistics/pdf/cancer\\_statistics\\_2025\\_fig\\_E.pdf](https://ganjoho.jp/public/qa_links/report/statistics/pdf/cancer_statistics_2025_fig_E.pdf).

### Disclosures

AH, AP, and DL are employees of PPD™ Evidera™ Real-World Data & Scientific Solutions, Thermo Fisher Scientific. KS is an employee of the General Incorporated Association Life Data Initiative. DW, SN, and TM are employees of NTT DATA Japan Corporation. MK, TS, and TK are employees of Clinical Study Support, Inc. Funding was provided by Thermo Fisher Scientific.