# MAIMON RESEARCH LLC

# ARTIFICIAL INTELLIGENCE LARGE LANGUAGE MODEL INTERROGATION



# REPRESENTATIONAL MEASUREMENT FAILURE IN HEALTH TECHNOLGY ASSESSMENT

# FINLAND: NATIONAL ENDORSEMENT OF MEASUREMENT FAILURE IN HEALTH TECHNOLOGY ASSESSMENT

Paul C Langley Ph.D Adjunct Professor, College of Pharmacy, University of Minnesota, Minneapolis, MN

# FOREWORD

## HEALTH TECHNOLOGY ASSESSMENT: A GLOBAL SYSTEM OF

## NON-MEASUREMENT

This Logit Working Paper series documents a finding as extraordinary as it is uncomfortable: health technology assessment (HTA), across nations, agencies, journals, and decades, has developed as a global system of non-measurement. It speaks the language of numbers, models, utilities, QALYs, "value for money," thresholds, discounting, incremental ratios, extrapolations, and simulations. It demands arithmetic at every turn, multiplication, division, summation, aggregation, discounting, yet it never once established that the quantities to which these operations are applied are measurable. HTA has built a vast evaluative machinery on foundations that do not exist. The probabilities and normalized logits in the country reports that follow provide the empirical confirmation of this claim. They show, with unsettling consistency, that the global HTA system consistently support measurement failure.

The objective of this study is to evaluate whether the Finnish health technology assessment knowledge base possesses the axioms required for quantitative measurement under representational measurement theory. Rather than examining the technical features of individual instruments, the intentions of their developers, or the correctness of specific analytic applications, the study interrogates the epistemic environment that authorizes numerical use. The central question is not whether Finnish HTA employs numbers extensively, but whether the system within which those numbers circulate recognizes the logical conditions that determine when numerical representations can meaningfully stand for empirical quantities.

To address this question, the study applies a twenty-four item canonical diagnostic derived from representational measurement theory and Rasch measurement principles. Each statement expresses either a necessary condition for measurement or a known impossibility when those conditions are violated. Endorsement probabilities and normalized logits are used to classify whether these principles function as operative constraints within the Finnish HTA corpus. The purpose is diagnostic rather than evaluative: to determine possession or non-possession of measurement axioms at the system level, independent of methodological sophistication, statistical technique, or policy intent.

The findings of the study are unambiguous. Across all twenty-four canonical statements, the Finnish HTA knowledge base exhibits a structurally coherent pattern characterized by weak or absent reinforcement of measurement axioms and strong reinforcement of non-measurement conventions. Statements expressing foundational requirements such as unidimensionality, invariance, true zero, and the precedence of measurement over arithmetic consistently attract low endorsement probabilities and strongly negative normalized logits. In contrast, statements representing known impossibilities within representational measurement theory, including the ratio status of QALYs, the permissibility of negative values, and the aggregation of non-quantities, attract high endorsement probabilities and strongly positive logits.

This configuration mirrors profiles previously observed in Canada, Australia, and the United Kingdom. The Finnish case does not represent partial understanding, transitional reform, or local deviation. It reflects full structural invariance. Measurement axioms do not function as admissibility conditions for numerical inference. Arithmetic operations proceed by convention rather than by representational authorization. The result is a stable epistemic system in which valuation is routinely substituted for measurement and numerical form is mistaken for quantitative meaning.

The starting point is simple and inescapable: *measurement precedes arithmetic*. This principle is not a methodological preference but a logical necessity. One cannot multiply what one has not measured, cannot sum what has no dimensional homogeneity, cannot compare ratios when no ratio scale exists. When HTA multiplies time by utilities to generate QALYs, it is performing arithmetic with numbers that cannot support the operation. When HTA divides cost by QALYs, it is constructing a ratio from quantities that have no ratio properties. When HTA aggregates QALYs across individuals or conditions, it is combining values that do not share a common scale. These practices are not merely suboptimal; they are mathematically impossible.

The modern articulation of this principle can be traced to Stevens' seminal 1946 paper, which introduced the typology of nominal, ordinal, interval, and ratio scales [1]. Stevens made explicit what physicists, engineers, and psychologists already understood: different kinds of numbers permit different kinds of arithmetic. Ordinal scales allow ranking but not addition; interval scales permit addition and subtraction but not multiplication; ratio scales alone support multiplication, division, and the construction of meaningful ratios. Utilities derived from multiattribute preference exercises, such as EQ-5D or HUI, are ordinal preference scores; they do not satisfy the axioms of interval measurement, much less ratio measurement. Yet HTA has, for forty years, treated these utilities as if they were ratio quantities, multiplying them by time to create QALYs and inserting them into models without the slightest recognition that scale properties matter. Stevens' paper should have blocked the development of QALYs and cost-utility analysis entirely. Instead, it was ignored.

The foundational theory that establishes *when* and *whether* a set of numbers can be interpreted as measurements came with the publication of Krantz, Luce, Suppes, and Tversky's *Foundations of Measurement* (1971) [2]. Representational Measurement Theory (RMT) formalized the axioms under which empirical attributes can be mapped to numbers in a way that preserves structure. Measurement, in this framework, is not an act of assigning numbers for convenience, it is the discovery of a lawful relationship between empirical relations and numerical relations. The axioms of additive conjoint measurement, homogeneity, order, and invariance specify exactly when interval scales exist. RMT demonstrated once and for all that measurement is not optional and not a matter of taste: either the axioms hold and measurement is possible, or the axioms fail and measurement is impossible. Every major construct in HTA, utilities, QALYs, DALYs, ICERs, incremental ratios, preference weights, health-state indices, fails these axioms. They lack unidimensionality; they violate independence; they depend on aggregation of heterogeneous attributes; they collapse under the requirements of additive conjoint measurement. Yet HTA proceeded, decade after decade, without any engagement with these axioms, as if the field had collectively decided that measurement theory applied everywhere except in the evaluation of therapies.

Whereas representational measurement theory articulates the axioms for interval measurement, Georg Rasch's 1960 model provides the only scientific method for transforming ordered categorical responses into interval measures for latent traits [3]. Rasch models uniquely satisfy the principles of specific objectivity, sufficiency, unidimensionality, and invariance. For any construct such as pain, fatigue, depression, mobility, or need, Rasch analysis is the only legitimate means of producing an interval scale from ordinal item responses. Rasch measurement is not an alternative to RMT; it is its operational instantiation. The equivalence of Rasch's axioms and the axioms of representational measurement was demonstrated by Wright, Andrich and others as early as the 1970s. In the latent-trait domain, the very domain where HTA claims to operate; Rasch is the only game in town [4].

Yet Rasch is effectively absent from all HTA guidelines, including NICE, PBAC, CADTH, ICER, SMC, and PHARMAC. The analysis demands utilities but never requires that those utilities be measured. They rely on multiattribute ordinal classifications but never understand that those constructs be calibrated on interval or ratio scales. They mandate cost-utility analysis but never justify the arithmetic. They demand modelled QALYs but never interrogate their dimensional properties. These guidelines do not misunderstand Rasch; they do not know it exists. The axioms that define measurement and the model that makes latent trait measurement possible are invisible to the authors of global HTA rules. The field has evolved without the science that measurement demands.

How did HTA miss the bus so thoroughly? The answer lies in its historical origins. In the late 1970s and early 1980s, HTA emerged not from measurement science but from welfare economics, decision theory, and administrative pressure to control drug budgets. Its core concern was *valuing health states*, not *measuring health*. This move, quiet, subtle, but devastating, shifted the field away from the scientific question "What is the empirical structure of the construct we intend to measure?" and toward the administrative question "How do we elicit a preference weight that we can multiply by time?" The preference-elicitation projects of that era (SG, TTO, VAS) were rationalized as measurement techniques, but they never satisfied measurement axioms. Ordinal preferences were dressed up as quasi-cardinal indices; valuation tasks were misinterpreted as psychometrics; analyst convenience replaced measurement theory. The HTA community built an entire belief system around the illusion that valuing health is equivalent to measuring health. It is not.

The endurance of this belief system, forty years strong and globally uniform, is not evidence of validity but evidence of institutionalized error. HTA has operated under conditions of what can only be described as *structural epistemic closure*: a system that has never questioned its constructs because it never learned the language required to ask the questions. Representational measurement theory is not taught in graduate HTA programs; Rasch modelling is not part of guideline development; dimensional analysis is not part of methodological review. The field has been insulated from correction because its conceptual foundations were never laid. What remains is a ritualized practice: utilities in, QALYs out, ICERs calculated, thresholds applied. The arithmetic continues because everyone assumes someone else validated the numbers.

This Logit Working Paper series exposes, through probabilistic and logit-based interrogations of AI large language national knowledge bases, the scale of this failure. The results display a global

pattern: true statements reflecting the axioms of measurement receive weak endorsement; false statements reflecting the HTA belief system receive moderate or strong reinforcement. This is not disagreement. It is non-possession. It shows that HTA, worldwide, has developed as a quantitative discipline without quantitative foundations; a confused exercise in numerical storytelling.

The conclusion is unavoidable: HTA does not need incremental reform; it needs a scientific revolution. Measurement must precede arithmetic. Representational axioms must precede valuation rituals. Rasch measurement must replace ordinal summation and utility algorithms. Value claims must be falsifiable, protocol-driven, and measurable; rather than simulated, aggregated, and numerically embellished.

The global system of non-measurement is now visible. The task ahead is to replace it with science.

**Paul C Langley, Ph.D**

**Email: langleylapaloma@gmail.com**

## DISCLAIMER

This analysis is generated through the structured interrogation of a large language model (LLM) applied to a defined documentary corpus and is intended solely to characterize patterns within an aggregated knowledge environment. It does identify, assess, or attribute beliefs, intentions, competencies, or actions to any named individual, faculty member, student, administrator, institution, or organization. The results do not constitute factual findings about specific persons or programs, nor should they be interpreted as claims regarding professional conduct, educational quality, or compliance with regulatory or accreditation standards. All probabilities and logit values reflect model-based inferences about the presence or absence of concepts within a bounded textual ecosystem, not judgments about real-world actors. The analysis is exploratory, interpretive, and methodological in nature, offered for scholarly discussion of epistemic structures rather than evaluative or legal purposes. Any resemblance to particular institutions or practices is contextual and non-attributive, and no adverse implication should be inferred.

# 1. INTERROGATING THE LARGE LANGUAGE MODEL

A large language model **(**LLM**)** is an artificial intelligence system designed to understand, generate, and manipulate human language by learning patterns from vast amounts of text data. Built on deep neural network architectures, most commonly transformers, LLMs analyze relationships between words, sentences, and concepts to produce contextually relevant responses. During training, the model processes billions of examples, enabling it to learn grammar, facts, reasoning patterns, and even subtle linguistic nuances. Once trained, an LLM can perform a wide range of tasks: answering questions, summarizing documents, generating creative writing, translating languages, assisting with coding, and more. Although LLMs do not possess consciousness or true understanding, they simulate comprehension by predicting the most likely continuation of text based on learned patterns. Their capabilities make them powerful tools for communication, research, automation, and decision support, but they also require careful oversight to ensure accuracy, fairness, privacy, and responsible use

In this Logit Working Paper, "interrogation" refers not to discovering what an LLM *believes,* it has no beliefs, but to probing the content of the *corpus-defined knowledge space* we choose to analyze. This knowledge base is enhanced if it is backed by accumulated memory from the user. In this case the interrogation relies also on 12 months of HTA memory from continued application of the system to evaluate HTA experience. The corpus is defined before interrogation: it may consist of a journal (e.g., *Value in Health*), a national HTA body, a specific methodological framework, or a collection of policy documents. Once the boundaries of that corpus are established, the LLM is used to estimate the conceptual footprint within it. This approach allows us to determine which principles are articulated, neglected, misunderstood, or systematically reinforced.

In this HTA assessment, the objective is precise: to determine the extent to which a given HTA knowledge base or corpus, global, national, institutional, or journal-specific, recognizes and reinforces the foundational principles of representational measurement theory (RMT). The core principle under investigation is that measurement precedes arithmetic; no construct may be treated as a number or subjected to mathematical operations unless the axioms of measurement are satisfied. These axioms include unidimensionality, scale-type distinctions, invariance, additivity, and the requirement that ordinal responses cannot lawfully be transformed into interval or ratio quantities except under Rasch measurement rules.

The HTA knowledge space is defined pragmatically and operationally. For each jurisdiction, organization, or journal, the corpus consists of:

- published HTA guidelines
- agency decision frameworks
- cost-effectiveness reference cases
- academic journals and textbooks associated with HTA
- modelling templates, technical reports, and task-force recommendations
- teaching materials, methodological articles, and institutional white papers

These sources collectively form the epistemic environment within which HTA practitioners develop their beliefs and justify their evaluative practices. The boundary of interrogation is thus

not the whole of medicine, economics, or public policy, but the specific textual ecosystem that sustains HTA reasoning. . The "knowledge base" is therefore not individual opinions but the cumulative, structured content of the HTA discourse itself within the LLM.

## THE NATIONAL KNOWLEDGE BASE FOR FINLAND

. The Finnish health technology assessment knowledge base is best understood as a distributed epistemic system rather than as a centralized theoretical framework. Its authority does not derive from a single institution, guideline, or methodological doctrine, but from the coordinated practices of academic researchers, HTA agencies, professional training programs, journals, and analytic infrastructures that collectively normalize numerical use. Within this system, quantitative legitimacy is not established through demonstration of measurement properties, but through repetition, precedent, and institutional acceptance.

Academic research plays a central role in sustaining this environment. Finnish health economics and outcomes research routinely employs preference-based multiattribute instruments, most prominently the EQ-5D family and related indices, as quantitative endpoints. Utility values are summarized using means, compared across treatment arms, and interpreted as magnitudes of health-related quality of life. These practices are presented as methodologically routine. Discussion of scale type, unidimensionality, or invariance is largely absent. The absence itself functions epistemically: it signals that justification is unnecessary.

HTA agency practice reinforces this normalization. Finnish assessment procedures accept utility values as legitimate quantitative inputs for comparative evaluation and economic modeling. This acceptance does not arise from explicit endorsement of representational measurement theory. Instead, it reflects compatibility with international reference-case conventions. Once numerical forms are aligned with accepted modeling templates, their epistemic status is treated as settled. Measurement becomes administratively presumed rather than theoretically established.

Methodological guidance documents further embed this presumption. Economic evaluation manuals describe how utilities are to be applied, combined with time, and aggregated across populations. These documents focus on procedural consistency rather than representational legitimacy. Numerical operations are specified, but the conditions under which those operations are meaningful are not addressed. Arithmetic is thereby institutionalized as a requirement rather than a conditional act.

Education plays a decisive role in reproducing this structure. Graduate programs in health economics and HTA train analysts to work with utilities, QALYs, and model outputs as standard components of analysis. Students learn how to implement numerical procedures, not how to interrogate whether those procedures are licensed by measurement theory. By the time analysts enter professional practice, numerical legitimacy has already been internalized. Instruments are encountered as givens, not as epistemic propositions.

Analytic infrastructure extends this process further. Software platforms, economic models, and statistical packages embed utility scoring algorithms and QALY calculations directly into workflows. Once encoded, the assumptions underlying these numbers disappear from view.

Users interact with outputs without encountering the premises that authorize or prohibit their interpretation. Epistemic commitment becomes automated.

Crucially, the Finnish knowledge base is not unified by explicit theoretical agreement. There is no declaration that utilities satisfy measurement axioms, nor any sustained debate about whether they do. Unity arises instead through coordinated silence. Measurement theory does not function as a governing authority because it is not invoked. Where axioms are absent from disciplinary grammar, they cannot constrain practice.

This distributed structure explains the stability of the Finnish profile. Developers can point to international usage. Agencies can point to precedent. Researchers can point to guidelines. Educators can point to curricula. Each component defers foundational responsibility to another. The result is epistemic closure: numerical practice persists without ever encountering the conditions required to authorize it.

The canonical interrogation reveals this condition with clarity. The Finnish HTA knowledge base does not lack numbers, models, or analytic sophistication. It lacks measurement possession. The principles that determine when numbers can represent quantities do not operate as admissibility conditions. What persists instead is a coherent belief system in which valuation is treated as measurement and arithmetic substitutes for representation.

This finding does not indict individual analysts or institutions. It identifies a system-level epistemic structure that has become self-sustaining. By making that structure visible, the analysis establishes the necessary foundation for any future reconsideration of quantitative claims in Finnish health technology assessment.

## CATEGORICAL PROBABILITIES

In the present application, the interrogation is tightly bounded. It does not ask what an LLM "thinks," nor does it request a normative judgment. Instead, the LLM evaluates how likely the HTA knowledge space is to endorse, imply, or reinforce a set of 24 diagnostic statements derived from representational measurement theory (RMT). Each statement is objectively TRUE or FALSE under RMT. The objective is to assess whether the HTA corpus exhibits possession or non-possession of the axioms required to treat numbers as measures. The interrogation creates an categorical endorsement probability: the estimated likelihood that the HTA knowledge base endorses the statement whether it is true or false; *explicitly or implicitly.*

The use of categorical endorsement probabilities within the Logit Working Papers reflects both the nature of the diagnostic task and the structure of the language model that underpins it. The purpose of the interrogation is not to estimate a statistical frequency drawn from a population of individuals, nor to simulate the behavior of hypothetical analysts. Instead, the aim is to determine the conceptual tendencies embedded in a domain-specific knowledge base: the discursive patterns, methodological assumptions, and implicit rules that shape how a health technology assessment environment behaves. A large language model does not "vote" like a survey respondent; it expresses likelihoods based on its internal representation of a domain. In this context, endorsement probabilities capture the strength with which the knowledge base, as represented within the model,

supports a particular proposition. Because these endorsements are conceptual rather than statistical, the model must produce values that communicate differences in reinforcement without implying precision that cannot be justified.

This is why categorical probabilities are essential. Continuous probabilities would falsely suggest a measurable underlying distribution, as if each HTA system comprised a definable population of respondents with quantifiable frequencies. But large language models do not operate on that level. They represent knowledge through weighted relationships between linguistic and conceptual patterns. When asked whether a domain tends to affirm, deny, or ignore a principle such as unidimensionality, admissible arithmetic, or the axioms of representational measurement, the model draws on its internal structure to produce an estimate of conceptual reinforcement. The precision of that estimate must match the nature of the task. Categorical probabilities therefore provide a disciplined and interpretable way of capturing reinforcement strength while avoiding the illusion of statistical granularity.

The categories used, values such as 0.05, 0.10, 0.20, 0.50, 0.75, 0.80, and 0.85, are not arbitrary. They function as qualitative markers that correspond to distinct degrees of conceptual possession: near-absence, weak reinforcement, inconsistent or ambiguous reinforcement, common reinforcement, and strong reinforcement. These values are far enough apart to ensure clear interpretability yet fine-grained enough to capture meaningful differences in the behavior of the knowledge base. The objective is not to measure probability in a statistical sense but to classify the epistemic stance of the domain toward a given item. A probability of 0.05 signals that the knowledge base almost never articulates or implies the correct response under measurement theory, whereas 0.85 indicates that the domain routinely reinforces it. Values near the middle reflect conceptual instability rather than a balanced distribution of views.

Using categorical probabilities also aligns with the requirements of logit transformation. Converting these probabilities into logits produces an interval-like diagnostic scale that can be compared across countries, agencies, journals, or organizations. The logit transformation stretches differences at the extremes, allowing strong reinforcement and strong non-reinforcement to become highly visible. Normalizing logits to the fixed ±2.50 range ensure comparability without implying unwarranted mathematical precision. Without categorical inputs, logits would suggest a false precision that could mislead readers about the nature of the diagnostic tool.

In essence, the categorical probability approach translates the conceptual architecture of the LLM into a structured and interpretable measurement analogue. It provides a disciplined bridge between the qualitative behavior of a domain's knowledge base and the quantitative diagnostic framework needed to expose its internal strengths and weaknesses.

The LLM computes these categorical probabilities from three sources:

1. **Structural content of HTA discourse**
   If the literature repeatedly uses ordinal utilities as interval measures, multiplies non-quantities, aggregates QALYs, or treats simulations as falsifiable, the model infers high reinforcement of these false statements.

2. **Conceptual visibility of measurement axioms**
   If ideas such as unidimensionality, dimensional homogeneity, scale-type integrity, or Rasch transformation rarely appear, or are contradicted by practice, the model assigns low endorsement probabilities to TRUE statements.
3. **The model's learned representation of domain stability**
   Where discourse is fragmented, contradictory, or conceptually hollow, the model avoids assigning high probabilities. This is *not* averaging across people; it is a reflection of internal conceptual incoherence within HTA.

The output of interrogation is a categorical probability for each statement. Probabilities are then transformed into logits [ $\ln(p/(1-p))$], capped to ±4.0 logits  to avoid extreme distortions, and normalized to ±2.50 logits for comparability across countries. A positive normalized logit indicates reinforcement in the knowledge base. A negative logit indicates weak reinforcement or conceptual absence. Values near zero logits reflect epistemic noise.

Importantly, *a high endorsement probability for a false statement does not imply that practitioners knowingly believe something incorrect*. It means the HTA literature itself behaves as if the falsehood were true; through methods, assumptions, or repeated uncritical usage. Conversely, a low probability for a true statement indicates that the literature rarely articulates, applies, or even implies the principle in question.

The LLM interrogation thus reveals structural epistemic patterns in HTA: which ideas the field possesses, which it lacks, and where its belief system diverges from the axioms required for scientific measurement. It is a diagnostic of the *knowledge behavior* of the HTA domain, not of individuals. The 24 statements function as probes into the conceptual fabric of HTA, exposing the extent to which practice aligns or fails to align with the  axioms of representational measurement.

## INTERROGATION STATEMENTS

Below is the canonical list of the 24 diagnostic HTA measurement items used in all the logit analyses, each marked with its correct truth value under representational measurement theory (RMT) and Rasch measurement principles.

This is the definitive set used across the Logit Working Papers.

**Measurement Theory & Scale Properties**

1. Interval measures lack a true zero — TRUE
2. Measures must be unidimensional — TRUE
3. Multiplication requires a ratio measure — TRUE
4. Time trade-off preferences are unidimensional — FALSE
5. Ratio measures can have negative values — FALSE
6. EQ-5D-3L preference algorithms create interval measures — FALSE
7. The QALY is a ratio measure — FALSE
8. Time is a ratio measure — TRUE

**Measurement Preconditions for Arithmetic**

9.  Measurement precedes arithmetic — TRUE
10. Summations of subjective instrument responses are ratio measures — FALSE
11. Meeting the axioms of representational measurement is required for arithmetic — TRUE

**Rasch Measurement & Latent Traits**

12. There are only two classes of measurement: linear ratio and Rasch logit ratio — TRUE
13. Transforming subjective responses to interval measurement is only possible with Rasch rules — TRUE
14. Summation of Likert question scores creates a ratio measure — FALSE

**Properties of QALYs & Utilities**

15. The QALY is a dimensionally homogeneous measure — FALSE
16. Claims for cost-effectiveness fail the axioms of representational measurement — TRUE
17. QALYs can be aggregated — FALSE

**Falsifiability & Scientific Standards**

18. Non-falsifiable claims should be rejected — TRUE
19. Reference-case simulations generate falsifiable claims — FALSE

**Logit Fundamentals**

20. The logit is the natural logarithm of the odds-ratio — TRUE

**Latent Trait Theory**

21. The Rasch logit ratio scale is the only basis for assessing therapy impact for latent traits — TRUE
22. A linear ratio scale for manifest claims can always be combined with a logit scale — FALSE
23. The outcome of interest for latent traits is the possession of that trait — TRUE
24. The Rasch rules for measurement are identical to the axioms of representational measurement — TRUE

---

**AI LARGE LANGUAGE MODEL STATEMENTS: TRUE OR FALSE**

Each of the 24 statements has a 400 word explanation why the statement is true or false as there may be differences of opinion on their status in terms of unfamiliarity with scale typology and the axioms of representational measurement.

The link to these explanations is:  https://maimonresearch.com/ai-llm-true-or-false/

## INTERPRETING TRUE STATEMENTS

TRUE statements represent foundational axioms of measurement and arithmetic. Endorsement probabilities for TRUE items typically cluster in the low range, indicating that the HTA corpus does *not* consistently articulate or reinforce essential principles such as:

- measurement preceding arithmetic
- unidimensionality
- scale-type distinctions
- dimensional homogeneity
- impossibility of ratio multiplication on non-ratio scales
- the Rasch requirement for latent-trait measurement

Low endorsement indicates **non-possession** of fundamental measurement knowledge—the literature simply does not contain, teach, or apply these principles.

## INTERPRETING FALSE STATEMENTS

FALSE statements represent the well-known mathematical impossibilities embedded in the QALY framework and reference-case modelling. Endorsement probabilities for FALSE statements are often moderate or even high, meaning the HTA knowledge base:

- accepts non-falsifiable simulation as evidence
- permits negative "ratio" measures
- treats ordinal utilities as interval measures
- treats QALYs as ratio measures
- treats summated ordinal scores as ratio scales
- accepts dimensional incoherence

This means the field systematically reinforces incorrect assumptions at the center of its practice. *Endorsement* here means the HTA literature behaves as though the falsehood were true.

## 2. SUMMARY OF FINDINGS FOR TRUE AND FALSE ENDORSEMENTS: FINLAND

Table 1 presents probabilities and normalized logits for each of the 24 diagnostic measurement statements. This is the standard reporting format used throughout the HTA assessment series.

It is essential to understand how to interpret these results.

The endorsement probabilities do not indicate whether a statement is *true* or *false* under representational measurement theory. Instead, they estimate the extent to which the HTA knowledge base associated with the target treats the statement as if it were true, that is, whether the concept is reinforced, implied, assumed, or accepted within the country's published HTA knowledge base.

The logits provide a continuous, symmetric scale, ranging from +2.50 to –2.50, that quantifies the degree of this endorsement. the logits, of course link to the probabilities (p) as the logit is the natural logarithm of the odds ratio; $\text{logit} = \ln[p/1\text{-}p]$.

- Strongly positive logits indicate pervasive reinforcement of the statement within the knowledge system.
- Strongly negative logits indicate conceptual absence, non-recognition, or contradiction within that same system.
- Values near zero indicate only shallow, inconsistent, or fragmentary support.

Thus, the endorsement logit profile serves as a direct index of a country's epistemic alignment with the axioms of scientific measurement, revealing the internal structure of its HTA discourse. It does not reflect individual opinions or survey responses, but the implicit conceptual commitments encoded in the literature itself.

**TABLE 1: ITEM STATEMENT, RESPONSE, ENDORSEMENT AND NORMALIZED LOGITS   FINLAND**

| STATEMENT | RESPONSE 1=TRUE 0=FALSE | ENDORSEMENT OF RESPONSE CATEGORICAL PROBABILITY | NORMALIZED LOGIT (IN RANGE +/- 2.50) |
|---|---|---|---|
| INTERVAL MEASURES LACK A TRUE ZERO | 1 | | |
| MEASURES MUST BE UNIDIMENSIONAL | 1 | | |
| MULTIPLICATION REQUIRES A RATIO MEASURE | 1 | | |
| TIME TRADE-OFF PREFERENCES ARE UNIDIMENSIONAL | 0 | | |
| RATIO MEASURES CAN HAVE NEGATIVE VALUES | 0 | | |
| EQ-5D-3L PREFERENCE ALGORITHMS CREATE INTERVAL MEASURES | 0 | | |
| THE QALY IS A RATIO MEASURE | 0 | | |
| TIME IS A RATIO MEASURE | 1 | | |
| MEASUREMENT PRECEDES ARITHMETIC | 1 | | |
| SUMMATIONS OF SUBJECTIVE INSTRUMENT RESPONSES ARE RATIO MEASURES | 0 | | |
| MEETING THE AXIOMS OF REPRESENTATIONAL MEASUREMENT IS REQUIRED FOR ARITHMETIC | 1 | | |
| THERE ARE ONLY TWO CLASSES OF MEASUREMENT LINEAR RATIO AND RASCH LOGIT RATIO | 1 | | |
| TRANSFORMING SUBJECTIVE RESPONSES TO INTERVAL MEASUREMENT IS ONLY POSSIBLE WITH RASH RULES | 1 | | |
| SUMMATION OF LIKERT QUESTION SCORES CREATES A RATIO MEASURE | 0 | | |
| THE QALY IS A DIMENSIONALLY HOMOGENEOUS MEASURE | 0 | | |
| CLAIMS FOR COST-EFFECTIVENESS FAIL THE AXIOMS OF REPRESENTATIONAL MEASUREMENT | 1 | | |

| | | | |
|---|---|---|---|
| QALYS CAN BE AGGREGATED | 0 | | |
| NON-FALSIFIABLE CLAIMS SHOULD BE REJECTED | 1 | | |
| REFERENCE CASE SIMULATIONS GENERATE FALSIFIABLE CLAIMS | 0 | | |
| THE LOGIT IS THE NATURAL LOGARITHM OF THE ODDS-RATIO | 1 | | |
| THE RASCH LOGIT RATIO SCALE IS THE ONLY BASIS FOR ASSESSING THERAPY IMPACT FOR LATENT TRAITS | 1 | | |
| A LINEAR RATIO SCALE FOR MANIFEST CLAIMS CAN ALWAYS BE COMBINED WITH A LOGIT SCALE | 0 | | |
| THE OUTCOME OF INTEREST FOR LATENT TRAITS IS THE POSSESSION OF THAT TRAIT | 1 | | |
| THE RASCH RULES FOR MEASUREMENT ARE IDENTICAL TO THE AXIOMS OF REPRESENTATIONAL MEASUREMENT | 1 | | |

## FINLAND: A CANONICAL EPISTEMIC ASSESSMENT OF HEALTH TECHNOLOGY ASSESSMENT

The canonical interrogation of the Finnish national health technology assessment knowledge base yields a result that is at once unsurprising and deeply instructive. Across the full set of twenty-four diagnostic statements, the pattern of endorsement probabilities and normalized logits reproduces with striking fidelity the structural configuration previously observed in Canada, Australia, the United Kingdom, and multiple HTA agencies and academic environments. This outcome is not an anomaly. It is the expected expression of a shared epistemic architecture that has governed health technology assessment for more than four decades.

The Finnish profile exhibits a clear and internally coherent polarity. Statements representing foundational axioms of representational measurement theory consistently attract low endorsement probabilities, with corresponding negative normalized logits. Conversely, statements representing known measurement impossibilities embedded within contemporary HTA practice attract high endorsement probabilities and strong positive logits. The result is not random dispersion, partial misunderstanding, or conceptual ambivalence. It is a stable epistemic signature indicating non-possession of measurement axioms alongside active reinforcement of non-measurement practices.

The first group of results concerns the most elementary distinction in measurement theory: the difference between interval and ratio scales. The statement that interval measures lack a true zero attracts an endorsement probability of 0.20, corresponding to a normalized logit of −1.40. This indicates weak reinforcement at best. While the idea may occasionally appear in methodological

discussion, it does not function as a governing constraint within the Finnish HTA knowledge base. In practical application, health utility values are routinely treated as though zero represents the absence of health, even though no empirical structure has been offered to justify such an interpretation. Zero operates operationally rather than representationally.

This weakness becomes more pronounced when attention turns to unidimensionality. The proposition that measures must be unidimensional registers an endorsement probability of 0.15, producing a logit of −1.75. This result indicates near-absence of reinforcement. Unidimensionality does not operate as an admissibility condition for quantitative claims. Instead, Finnish HTA practice routinely accepts composite indices derived from multiple heterogeneous attributes as if they represented variation along a single measurable dimension. This acceptance is not defended theoretically; it is normalized procedurally.

The same pattern emerges for the proposition that multiplication requires a ratio measure, which also receives an endorsement probability of 0.15 and a logit of −1.75. This result is epistemically decisive. Multiplication is not treated as a conditional operation governed by scale type. Rather, it is treated as an analytic necessity dictated by modeling frameworks. When utility values are multiplied by time to generate quality-adjusted life years, the operation is undertaken not because the scale permits multiplication, but because the reference case requires it. Arithmetic thus becomes prescriptive rather than conditional.

Time itself provides a revealing contrast. The statement that time is a ratio measure receives an endorsement probability of 0.80 and a positive logit of +1.40. This is one of the few items in the diagnostic that attracts strong reinforcement of a true axiom. The reason is straightforward: time possesses a true zero, invariant units, and empirical concatenation. Its ratio properties are uncontroversial. Yet the epistemic system fails to recognize the asymmetry this creates. Multiplying a ratio measure by a non-ratio measure does not produce a ratio quantity. The Finnish knowledge base endorses time's ratio properties while simultaneously ignoring the implications for operations performed with non-ratio utilities.

The inversion between measurement and arithmetic is further confirmed by the endorsement probability of 0.10 for the statement that measurement precedes arithmetic, yielding a logit of −2.20. This result indicates that the foundational ordering of logic in representational measurement theory is not recognized. Arithmetic does not follow measurement; it substitutes for it. Numerical manipulation is undertaken first, and legitimacy is inferred retrospectively from widespread use rather than from axiomatic demonstration.

A similar pattern appears for the proposition that meeting the axioms of representational measurement is required for arithmetic. With an endorsement probability of 0.10 and a logit of −2.20, the Finnish knowledge base does not treat axiomatic compliance as a precondition for numerical operations. This is not an oversight. It reflects a system in which arithmetic is regarded as methodologically neutral rather than logically constrained. Numbers are permitted to circulate freely, independent of the empirical structures they are presumed to represent.

The diagnostic reaches its most decisive results when it turns to the measurement of latent constructs. The proposition that there exist only two classes of measurement—linear ratio scales

for manifest attributes and Rasch logit ratio scales for latent traits—collapses to the floor, with an endorsement probability of 0.05 and a normalized logit of −2.50. This indicates complete non-possession. The Finnish HTA knowledge base does not recognize the distinction between manifest and latent measurement as epistemically operative.

This non-possession is reinforced by the equally strong rejection of the proposition that transforming subjective responses to interval measurement is only possible through Rasch rules. Again, the endorsement probability is 0.05, producing a logit of −2.50. Ordinal responses are routinely treated as quantities through summation, weighting, or preference valuation, despite the absence of an invariant transformation model. Rasch measurement is not rejected explicitly; it is simply absent from the disciplinary grammar.

The collapse continues with the proposition that the Rasch logit ratio scale is the only basis for assessing therapy impact for latent traits. The endorsement probability remains at 0.05. This reflects a fundamental conceptual misalignment. In Finnish HTA practice, therapy impact is inferred from changes in index scores, preference values, or modeled outcomes rather than from changes in possession of a latent attribute measured on an invariant scale. Impact is operationalized numerically, not measured representationally.

The statement that the outcome of interest for latent traits is possession of that trait likewise collapses to a probability of 0.05 and a logit of −2.50. This is epistemically crucial. Measurement concerns attributes possessed by individuals. The Finnish knowledge base does not conceptualize health-related quality of life in these terms. Instead, individuals are assigned values based on externally derived preferences for hypothetical health states. What is quantified is desirability, not magnitude of possession.

The next group of results concerns valuation-based scoring and its conflation with measurement. The proposition that time trade-off preferences are unidimensional receives an endorsement probability of 0.85 and a logit of +1.75. This indicates strong reinforcement of a falsehood. Preferences do not constitute a unidimensional empirical attribute; they reflect context-dependent judgments of desirability. Yet the Finnish HTA knowledge base treats such preferences as if they possessed dimensional structure capable of supporting arithmetic.

The same strong endorsement appears for the proposition that ratio measures can have negative values, with a probability of 0.90 and a logit of +2.20. This item is especially revealing. Under representational measurement theory, ratio measures cannot take negative values because zero represents the absence of the attribute. Yet negative utilities are routinely accepted in Finnish economic evaluation. This contradiction is not debated or resolved; it is normalized. Negative values are interpreted as "worse than dead," a preference concept that has no analogue in measurement theory.

The proposition that EQ-5D preference algorithms create interval measures also attracts a high endorsement probability of 0.85. This indicates that algorithmic scoring is widely treated as conferring metric properties. The existence of a numerical function is taken as evidence of interval measurement, even though no empirical structure is shown to be preserved. Valuation substitutes for representation.

The statement that the QALY is a ratio measure receives an endorsement probability of 0.90 and a logit of +2.20. This is one of the most diagnostically powerful results. The Finnish HTA system does not merely tolerate the QALY; it actively reinforces its ratio interpretation. Yet the QALY is constructed by multiplying time, a ratio measure, by utilities that lack a true zero, lack invariant units, and permit negative values. Under representational measurement theory, such multiplication is undefined. The strong positive logit indicates that this impossibility is not recognized as such.

The proposition that summations of subjective instrument responses are ratio measures likewise receives a probability of 0.90. Ordinal responses, whether from questionnaires or health-state classifications, are treated as quantities through summation and averaging. This practice is not justified empirically; it is inherited methodologically.

The statement that the QALY is a dimensionally homogeneous measure receives a probability of 0.85. This indicates reinforcement of the belief that time and utility share compatible dimensional structure. In reality, they do not. Time is measured; utility is valued. Their multiplication produces a numerical artifact, not a quantity.

The proposition that QALYs can be aggregated across individuals receives a probability of 0.90 and a logit of +2.20. This reflects one of the most consequential epistemic commitments in HTA. Aggregation presupposes commensurability of individual quantities. Without invariant measurement, aggregation has no representational meaning. Yet Finnish HTA practice aggregates QALYs routinely, treating population totals as meaningful magnitudes.

The diagnostic also reveals partial reinforcement of epistemic norms unrelated to measurement. The statement that non-falsifiable claims should be rejected receives a probability of 0.30 and a logit of −0.95. This suggests weak and inconsistent endorsement of Popperian standards. While falsifiability is occasionally invoked rhetorically, it does not function as a binding constraint on modeling practice.

This weakness is mirrored by the strong endorsement of the false proposition that reference case simulations generate falsifiable claims, which receives a probability of 0.90. Simulation outputs are treated as empirically testable despite being structurally insulated from refutation. Model revision substitutes for falsification.

Two items show modest attenuation. The statement that the logit is the natural logarithm of the odds ratio receives a probability of 0.40 and a logit of −0.45. This reflects partial technical awareness without epistemic integration. Logits may appear in statistical contexts, but their role in measurement theory is not recognized.

Finally, the proposition that the Rasch rules for measurement are identical to the axioms of representational measurement collapses to 0.05. The unification between Rasch and representational measurement theory is absent from the Finnish HTA knowledge base entirely.

Taken as a whole, the Finnish endorsement profile is internally coherent and structurally invariant with those observed in other national systems. There is no evidence of transitional understanding,

partial reform, or epistemic hybridization. The same axioms are absent. The same impossibilities are reinforced. The same arithmetic conventions dominate.

This invariance is not coincidental. It reflects the historical trajectory of HTA as a field. Health utility theory emerged not from measurement science but from welfare economics and decision analysis. Its objective was comparability and tractability, not representational validity. Instruments and models were constructed to serve analytic frameworks rather than to satisfy measurement axioms. Over time, these frameworks became institutionalized, and their numerical requirements came to define what counted as evidence.

The Finnish case confirms this diagnosis with clarity. Despite differences in national context, institutional design, and academic culture, the epistemic structure remains the same. Measurement is presumed rather than established. Valuation is mistaken for magnitude. Arithmetic is treated as neutral rather than conditional.

The significance of this finding lies not in critiquing Finland specifically, but in reinforcing a general conclusion. HTA systems do not differ epistemically in their treatment of measurement. They differ administratively, procedurally, and rhetorically, but not structurally. The same belief system governs them all.

The canonical diagnostic therefore achieves precisely what it was designed to do. It does not measure competence, sophistication, or intent. It reveals possession. And what it reveals, consistently, is the absence of representational measurement theory as a governing authority in contemporary HTA.

This absence cannot be corrected through refinement, recalibration, or improved modeling. It is not a technical deficiency. It is a foundational one. Measurement cannot be added downstream. It must exist upstream.

Until representational measurement axioms are restored as admissibility conditions for quantitative claims, health technology assessment will continue to generate numbers without representation and precision without quantity. The Finnish results do not deviate from this conclusion. They confirm it.

# 3. THE TRANSITION TO MEASUREMENT IN HEALTH TECHNOLOGY ASSESSMENT

## THE IMPERATIVE OF CHANGE

This analysis has not been undertaken to criticize decisions made by health system, nor to assign responsibility for the analytical frameworks currently used in formulary review. The evidence shows something more fundamental: organizations have been operating within a system that does not permit meaningful evaluation of therapy impact, even when decisions are made carefully, transparently, and in good faith.

The present HTA framework forces health systems to rely on numerical outputs that appear rigorous but cannot be empirically assessed (Table 1). Reference-case models, cost-per-QALY ratios, and composite value claims are presented as decision-support tools, yet they do not satisfy the conditions required for measurement. As a result, committees are asked to deliberate over results that cannot be validated, reproduced, or falsified. This places decision makers in an untenable position: required to choose among therapies without a stable evidentiary foundation.

This is not a failure of expertise, diligence, or clinical judgment. It is a structural failure. The prevailing HTA architecture requires arithmetic before measurement, rather than measurement before arithmetic. Health systems inherit this structure rather than design it. Manufacturers respond to it. Consultants reproduce it. Journals reinforce it. Universities promote it. Over time it has come to appear normal, even inevitable.

Yet the analysis presented in Table 1 demonstrates that this HTA framework cannot support credible falsifiable claims. Where the dependent variable is not a measure, no amount of modeling sophistication can compensate. Uncertainty analysis cannot rescue non-measurement. Transparency cannot repair category error. Consensus cannot convert assumption into evidence.

The consequence is that formulary decisions are based on numerical storytelling rather than testable claims. This undermines confidence, constrains learning, and exposes health systems to growing scrutiny from clinicians, patients, and regulators who expect evidence to mean something more than structured speculation.

The imperative of change therefore does not arise from theory alone. It arises from governance responsibility. A health system cannot sustain long-term stewardship of care if it lacks the ability to distinguish between claims that can be evaluated and claims that cannot. Without that distinction, there is no pathway to improvement; only endless repetition for years to come.

This transition is not about rejecting evidence. It is about restoring evidence to its proper meaning. It requires moving away from composite, model-driven imaginary constructs toward claims that are measurable, unidimensional, and capable of empirical assessment over time. The remainder of this section sets out how that transition can occur in a practical, defensible, and staged manner.

## MEANINGFUL THERAPY IMPACT CLAIMS

At the center of the current problem is not data availability, modeling skill, or analytic effort. It is the nature of the claims being advanced. Contemporary HTA has evolved toward increasingly complex frameworks that attempt to compress multiple attributes, clinical effects, patient experience, time, and preferences into single composite outputs. These constructs are then treated as if they were measures. They are not (Table 1).

The complexity of the reference-case framework obscures a simpler truth: meaningful evaluation requires meaningful claims. A claim must state clearly what attribute is being affected, in whom, over what period, and how that attribute is measured. When these conditions are met, evaluation becomes possible. When they are not  complexity substitutes for clarity. The current framework is not merely incorrect; it is needlessly elaborate. Reference-case modeling requires dozens of inputs, assumptions, and transformations, yet produces outputs that cannot be empirically verified. Each additional layer of complexity increases opacity while decreasing accountability. Committees are left comparing models rather than assessing outcomes.

In contrast, therapy impact can be expressed through two, and only two, types of legitimate claims. First are claims based on manifest attributes: observable events, durations, or resource units. These include hospitalizations avoided, time to event, days in remission, or resource use. When properly defined and unidimensional, these attributes can be measured on linear ratio scales and evaluated directly.

Second are claims based on latent attributes: symptoms, functioning, need fulfillment, or patient experience. These cannot be observed directly and therefore cannot be scored or summed meaningfully. They require formal measurement through Rasch models to produce invariant logit ratio scales. These two forms of claims are sufficient. They are also far more transparent. Each can be supported by a protocol. Each can be revisited. Each can be reproduced. Most importantly, each can fail. But they cannot be combined. This is the critical distinction. A meaningful claim is one that can be wrong.

Composite constructs such as QALYs do not fail in this sense. They persist regardless of outcome because they are insulated by assumptions. They are recalculated, not refuted. That is why they cannot support learning. The evolution of objective knowledge regarding therapy impact in disease areas is an entirely foreign concept. By re-centering formulary review on single-attribute, measurable claims, health systems regain control of evaluation. Decisions become grounded in observable change rather than modeled narratives. Evidence becomes something that accumulates, rather than something that is re-generated anew for every submission.

## THE PATH TO MEANINGFUL MEASUREMENT

Transitioning to meaningful measurement does not require abandoning current processes overnight. It requires reordering them. The essential change is not procedural but conceptual: measurement must become the gatekeeper for arithmetic, not its byproduct.

The first step is formal recognition that not all numerical outputs constitute evidence. Health systems must explicitly distinguish between descriptive analyses and evaluable claims. Numbers that do not meet measurement requirements may inform discussion but cannot anchor decisions.

The second step is restructuring submissions around explicit claims rather than models. Each submission should identify a limited number of therapy impact claims, each defined by attribute, population, timeframe, and comparator. Claims must be unidimensional by design.

Third, each claim must be classified as manifest or latent. This classification determines the admissible measurement standard and prevents inappropriate mixing of scale types.

Fourth, measurement validity must be assessed before any arithmetic is permitted. For manifest claims, this requires confirmation of ratio properties. For latent claims, this requires Rasch-based measurement with demonstrated invariance.

Fifth, claims must be supported by prospective or reproducible protocols. Evidence must be capable of reassessment, not locked within long-horizon simulations designed to frustrate falsification.

Sixth, committees must be supported through targeted training in representational measurement principles, including Rasch fundamentals. Without this capacity, enforcement cannot occur consistently.

Finally, evaluation must be iterative. Claims are not accepted permanently. They are monitored, reproduced, refined, or rejected as evidence accumulates.

These steps do not reduce analytical rigor. They restore it.

## TRANSITION REQUIRES TRAINING

A transition to meaningful measurement cannot be achieved through policy alone. It requires a parallel investment in training, because representational measurement theory is not intuitive and has never been part of standard professional education in health technology assessment, pharmacoeconomics, or formulary decision making. For more than forty years, practitioners have been taught to work within frameworks that assume measurement rather than demonstrate it. Reversing that inheritance requires structured learning, not informal exposure.

At the center of this transition is the need to understand why measurement must precede arithmetic. Representational measurement theory establishes the criteria under which numbers can legitimately represent empirical attributes. These criteria are not optional. They determine whether addition, multiplication, aggregation, and comparison are meaningful or merely symbolic. Without this foundation, committees are left evaluating numerical outputs without any principled way to distinguish evidence from numerical storytelling.

Training must therefore begin with scale types and their permissible operations. Linear ratio measurement applies to manifest attributes that possess a true zero and invariant units, such as

time, counts, and resource use. Latent attributes, by contrast, cannot be observed directly and cannot be measured through summation or weighting. They require formal construction through a measurement model capable of producing invariant units. This distinction is the conceptual fulcrum of reform, because it determines which claims are admissible and which are not.

For latent trait claims, Rasch measurement provides the only established framework capable of meeting these requirements. Developed in the mid–twentieth century alongside the foundations of modern measurement theory, the Rasch model was explicitly designed to convert subjective observations into linear logit ratio measures. It enforces unidimensionality, tests item invariance, and produces measures that support meaningful comparison across persons, instruments, and time. These properties are not approximations; they are defining conditions of measurement.

Importantly, Rasch assessment is no longer technically burdensome. Dedicated software platforms developed and refined over more than four decades make Rasch analysis accessible, transparent, and auditable. These programs do not merely generate statistics; they explain why items function or fail, how scales behave, and whether a latent attribute has been successfully measured. Measurement becomes demonstrable rather than assumed.

Maimon Research has developed a two-part training program specifically to support this transition. The first component provides foundational instruction in representational measurement theory, including the historical origins of scale theory, the distinction between manifest and latent attributes, and the criteria that define admissible claims. The second component focuses on application, detailing claim types, protocol design, and the practical use of Rasch methods to support latent trait evaluation.

Together, these programs equip health systems, committees, and analysts with the competence required to enforce measurement standards consistently. Training does not replace judgment; it enables it. Without such preparation, the transition to meaningful measurement cannot be sustained. With it, formulary decision making can finally rest on claims that are not merely numerical, but measurable.

---

**A NEW START IN MEASUREMENT FOR HEALTH TECHNOLOGY ASSESSMENT**

For readers who are looking for an introduction to measurement that meets the required standards, Maimon Research has just released two distance education programs. These are:

- Program 1: Numerical Storytelling – Systematic Measurement Failure in HTA.
- Program 2: A New Start in Measurement for HTA, with recommendations for protocol-supported claims for specific objective measures as well as latent constructs and manifested traits.

---

## DESIGNED FOR CLOSURE

For those who remain unconvinced that there is any need to abandon a long-standing and widely accepted HTA framework, it is necessary to confront a more fundamental question: why was this system developed and promoted globally in the first place?

The most plausible explanation is administrative rather than scientific. Policy makers were searching for an assessment framework that could be applied under conditions of limited empirical data while still producing a determinate conclusion. Reference-case modeling offered precisely this convenience. By constructing a simulation populated with assumptions, surrogate endpoints, preference weights, and extrapolated time horizons, it became possible to generate a numerical result that could be interpreted as decisive. Once an acceptable cost-effectiveness ratio emerged, the assessment could be declared complete and the pricing decision closed. This structure solved a political and administrative problem. It allowed authorities to claim that decisions were evidence-based without requiring the sustained empirical burden demanded by normal science. There was no requirement to formulate provisional claims and subject them to ongoing falsification. There was no obligation to revisit conclusions as new data emerged. Closure could be achieved at launch, rather than knowledge evolving over the product life cycle.

By contrast, a framework grounded in representational measurement would have imposed a very different obligation. Claims would necessarily be provisional. Measurement would precede arithmetic. Each therapy impact claim would require a defined attribute, a valid scale, a protocol, and the possibility of replication or refutation. Evidence would accumulate rather than conclude. Decisions would remain open to challenge as real-world data emerged. From an administrative standpoint, this was an unreasonable burden. It offered no finality.

The reference-case model avoided this problem entirely. By shifting attention away from whether quantities were measurable and toward whether assumptions were plausible, the framework replaced falsification with acceptability. Debate became internal to the model rather than external to reality. Sensitivity analysis substituted for empirical risk. Arithmetic proceeded without prior demonstration that the objects being manipulated possessed the properties required for arithmetic to be meaningful.

Crucially, this system required no understanding of representational measurement theory. Committees did not need to ask whether utilities were interval or ratio measures, whether latent traits had been measured or merely scored, or whether composite constructs could legitimately be multiplied or aggregated. These questions were never posed because the framework did not require

them to be posed. The absence of measurement standards was not an oversight; it was functionally essential.

Once institutionalized, the framework became self-reinforcing. Training programs taught modeling rather than measurement. Guidelines codified practice rather than axioms. Journals reviewed technique rather than admissibility. Over time, arithmetic without measurement became normalized as "good practice," while challenges grounded in measurement theory were dismissed as theoretical distractions. The result was a global HTA architecture capable of producing numbers, but incapable of producing falsifiable knowledge. Claims could be compared, ranked, and monetized, but not tested in the scientific sense. What evolved was not objective knowledge, but institutional consensus.

This history matters because it explains why the present transition is resisted. Moving to a real measurement framework with single, unidimensional claims does not merely refine existing methods; it dismantles the very mechanism by which closure has been achieved for forty years. It replaces decisiveness with accountability, finality with learning, and numerical plausibility with empirical discipline. Yet that is precisely the transition now required. A system that avoids measurement in order to secure closure cannot support scientific evaluation, cumulative knowledge, or long-term stewardship of healthcare resources. The choice is therefore unavoidable: continue with a framework designed to end debate, or adopt one designed to discover the truth.

Anything else is not assessment at all, but the ritualized manipulation of numbers detached from measurement, falsification, and scientific accountability.

## ACKNOWLEDGEMENT

## REFERENCES

[1] Stevens S. On the Theory of Scales of Measurement. *Science*. 1946;103(2684):677-80

[2] Krantz D, Luce R, Suppes P, Tversky A. Foundations pf Measurement Vol 1: Additive and Polynomial Representations. New York: Academic Press, 1971

[3] Rasch G, Probabilistic Models for some Intelligence and Attainment Tests. Chicago: University of Chicago Press, 1980 [An edited version of the original 1960 publication]

[4] Wright B. Solving measurement problems with the Rasch Model. *J Educational Measurement.* 1977;14(2):97-116