

Training-Free Adaptation of New-Generation LLMs using Legacy Clinical Models

Sasha Ronaghi
Stanford University

SRONAGHI@STANFORD.EDU

Chloe Stanwyck
Stanford University

CHLOEO@STANFORD.EDU

Asad Aali
Stanford University

ASADAALI@STANFORD.EDU

Amir Ronaghi
MemorialCare

ARONAGHI1@GMAIL.COM

Miguel Angel Fuentes Hernandez
Stanford University

MIGUFUEN@STANFORD.EDU

Tina Hernandez Boussard
Stanford University

BOUSSARD@STANFORD.EDU

Emily Alsentzer
Stanford University

ELSANTZER@STANFORD.EDU

Abstract

Adapting language models to the clinical domain through continued pretraining and instruction tuning requires costly retraining for each new model generation. We propose *Cross-Architecture Proxy Tuning* (CAPT), a model-ensembling approach that enables training-free adaptation of state-of-the-art general-domain models using existing clinical models. CAPT supports models with disjoint vocabularies, leveraging contrastive decoding to selectively inject clinically relevant signals while preserving the general-domain model’s reasoning and fluency. On six clinical classification and text-generation tasks, CAPT with a new-generation general-domain model and an older-generation clinical model consistently outperforms both models individually and state-of-the-art ensembling approaches (average +17.6% over UniTE, +41.4% over proxy tuning across tasks). Through token-level analysis and physician case studies, we demonstrate that CAPT amplifies clinically actionable language, reduces context errors, and increases clinical specificity. This technique especially benefits healthcare institutions with constrained computational capacity that cannot support iterative clinical training and want to adopt emerging general-domain model advances.

Data and Code Availability All datasets used in this paper are publicly available and shown in Table 1. Code available at: https://github.com/sronaghi/training_free_adaptation.

Institutional Review Board (IRB) This research uses publicly available datasets and does not require IRB approval.

1. Introduction

Despite advances in general-domain language models, their application to clinical settings remains limited by hallucinations, omission of critical details, and failures in clinical reasoning (Lehman et al., 2023; Hager et al., 2024; Asgari et al., 2025). These shortcomings arise because pretraining corpora contain limited representations of clinical data, such as electronic health records, due to privacy constraints (Singhal et al., 2023a). Moreover, reliance on large-scale Internet text during training encodes biases, outdated information, or incorrect medical knowledge (Alber et al., 2025; Zack et al., 2024; Wu et al., 2025a).

Domain adaptation techniques such as continued pretraining and instruction tuning address these challenges but require resource-intensive re-training for each new model generation (e.g., MedPalm 1 → Med-

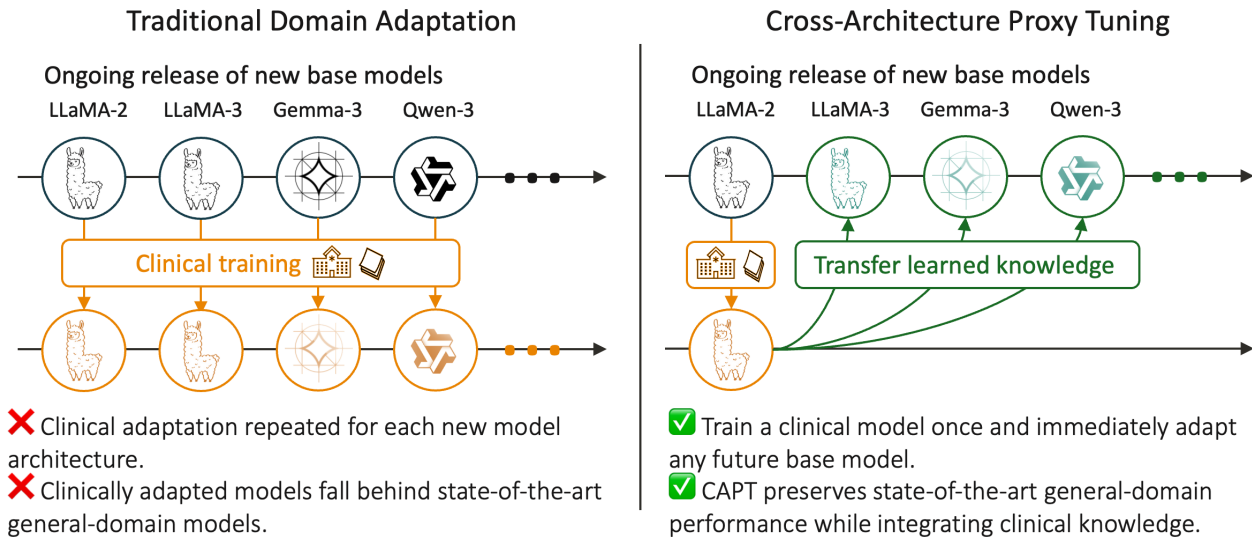


Figure 1: **Cross-Architecture Proxy Tuning (CAPT) Overview:** Conventional approaches require costly retraining for each new architecture (left), causing clinical models to lag behind general-domain counterparts. CAPT transfers knowledge from a single clinical model to new-generation models without additional training (right).

Palm 2) (Singhal et al., 2023b). This creates a lag between advances in base model capabilities and their clinical applicability, with substantial computational resources reinvested in domain adaptation for each new architecture. As a result, on comprehensive clinical benchmarks, general-domain state-of-the-art models often outperform clinically adapted models (Wu et al., 2025b).

We explore model ensembling as a training-free approach to combine the advanced reasoning of new-generation general-domain models with the domain knowledge of legacy clinical models. Prior works show that integrating decoding-time probability distributions of general-domain models with heterogeneous architectures yields efficient performance gains (Yao et al., 2025). However, these approaches assume largely overlapping model capabilities and have not been explored for combining models with distinct strengths, where information must be selectively integrated rather than uniformly aggregated. For example, we seek to incorporate the learned clinical knowledge of an older-generation clinical model without inheriting its limitations in instruction following, reasoning, or robustness due to its architecture or catastrophic forgetting (Kirkpatrick et al., 2017).

Contrastive decoding offers a mechanism for isolating a model’s strengths by choosing tokens that

maximize the likelihood difference between expert and amateur models, amplifying the expert’s behavior (Li et al., 2023). Proxy tuning leverages contrastive decoding to efficiently tune a large, pre-trained model by combining its logit distribution with the distributional delta between a smaller fine-tuned model and its untuned, base counterpart (Liu et al., 2024). However, proxy tuning requires shared tokenization across models, restricting architectural diversity. This constraint prevents reuse of domain-adapted models built on older architectures, restricting adoption in resource-limited settings.

Here, we propose **Cross-Architecture Proxy Tuning (CAPT)**, a probability-level ensembling method that supports models with disjoint vocabularies and leverages contrastive decoding to selectively incorporate the specialized knowledge of a domain-adapted model (Figure 1). Our contributions are three-fold:

- We introduce CAPT which enables reuse of legacy domain-adapted clinical models for training-free adaptation of newer-generation models. On six clinical classification and text-generation benchmarks, CAPT consistently outperforms prior ensembling approaches, with an

average improvement of 17.6% over UniTE and 41.4% over proxy tuning across metrics.

- Through a token-level analysis, we illustrate that CAPT selectively integrates the clinical model’s knowledge for tokens related to clinical decision-making and documentation style, while the general model controls linguistic structure and formatting tokens.
- In case studies by two board-certified physicians, we demonstrate that CAPT-generated outputs contain more precise clinical terminology, context-appropriate recommendations, and improved clinical accuracy.

2. Related Works

Clinical Domain Adaptation. Clinical models such as GatorTron, clinicalBERT, NYUTron, MeLLaMA, and Clinical-T5 have demonstrated the effectiveness of continued pretraining on large volumes of unlabeled medical text (Yang et al., 2022; Alsentzer et al., 2019; Jiang et al., 2023; Xie et al., 2024; Lehman and Johnson, 2023). Following continued pretraining, supervised fine-tuning can further improve performance on specific clinical tasks, instruction following, and alignment with human preferences (Han et al., 2025; Zhang et al., 2025; Singhal et al., 2023a; Wang et al., 2025; Zhang et al., 2023). While effective, these approaches are highly resource-intensive, especially continued pre-training; for example, MeLLaMA conducted continued pretraining on LLaMA-2 base models using 129B tokens and 160×80 GB A100 GPUs (Xie et al., 2024). Our work complements these techniques by addressing their limited transferability, as these costly adaptations must be repeated for each new base model.

Model Ensembling. Probability-level model ensembling methods for heterogeneous architectures have primarily been explored with multiple general-domain models. Existing works focus on optimizing efficiency and vocabulary alignment: EVA learns mappings between LLM vocabularies through their existing vocabulary overlap (Xu et al., 2024), and DeePEn projects the probability distributions of multiple models onto a unified space (Huang et al., 2024). UNiTE, the strongest-performing method to date, unions each model’s top- k tokens at every decoding step and combines probabilities via re-tokenization (Yao et al., 2025). Prior analyses show ensembling

methods only outperform the ensemble’s strongest model when performance gaps are within 10% (Yao et al., 2025). This limits applicability when combining new- and old-generation architectures where gaps can arise from architectural advances alone (Bedi et al., 2025). In contrast, CAPT is designed for models with asymmetric capabilities and enables selective integration of learned information from older domain models.

Contrastive Decoding. Contrastive decoding selects tokens that are more likely under a strong language model by contrasting its logit or probability output distribution against that of a weaker model. Prior work operationalizes this idea by combining the logits of a large model with the distributional delta between a smaller strong model and its untuned base counterpart, steering generation toward desirable behaviors such as reduced toxicity (Liu et al., 2021), increased helpfulness (Mitchell et al., 2023), selective unlearning (Suriyakumar et al., 2025), and improved coding, question answering, and mathematical reasoning (Liu et al., 2024). These approaches typically assume shared vocabularies and are evaluated in task- or instruction-tuning regimes that remain in-distribution with respect to the large model’s pre-training data. In contrast, our approach supports models with disjoint vocabularies and is evaluated on clinical domain adaptation, a substantially out-of-distribution setting (Kim et al., 2025).

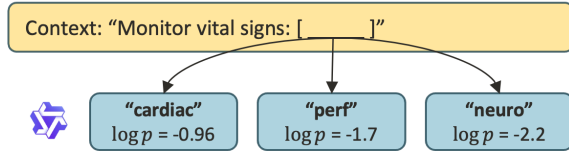
3. Methods

3.1. Cross-Architecture Proxy Tuning

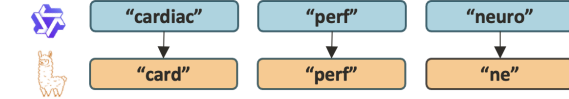
Overview. We hypothesize that general-domain models encode substantial medical knowledge but lack exposure to clinical practice patterns and stylistic conventions reflected in clinical notes. Accordingly, the new-generation general-domain model should lead generation to preserve fluency and reasoning, while the clinical model selectively interjects for tokens associated with domain-specific reasoning or stylistic patterns. To achieve this, CAPT re-ranks the top- k candidate tokens proposed by the general-domain model using a log-probability offset defined by the difference between the clinically trained model and its untrained base counterpart.

At each decoding step, we select the top- k next-token candidates from the new-generation general-domain model. Each candidate token is re-tokenized

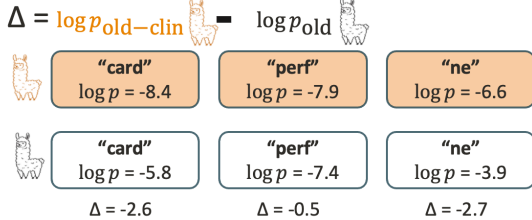
1. New model proposes Top- k next token candidates



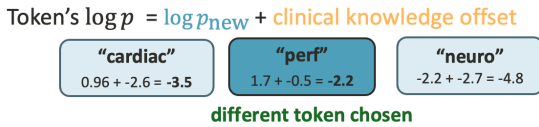
2. Map to clinical model’s vocabulary



3. Calculate clinical knowledge offset



4. Re-rank and select best token



Formal definition. Let M_{new} be a new-generation general-domain language model with tokenizer \mathcal{T}_{new} and vocabulary \mathcal{V}_{new} . Let $M_{\text{old-clin}}$ denote an older-generation clinically trained model and M_{old} its untrained base counterpart; both share tokenizer \mathcal{T}_{old} and vocabulary \mathcal{V}_{old} .

Given a context $\mathbf{x}_{1:t}$, the models define next-token log-probabilities $\log p_{\text{new}}(\cdot | \mathbf{x}_{1:t})$ over \mathcal{V}_{new} and $\log p_{\text{old-clin}}(\cdot | \mathbf{x}_{1:t})$, $\log p_{\text{old}}(\cdot | \mathbf{x}_{1:t})$ over \mathcal{V}_{old} . We restrict computation to the candidate set $\mathcal{C}_t = \text{Top-}k(\log p_{\text{new}}(\cdot | \mathbf{x}_{1:t})) \subseteq \mathcal{V}_{\text{new}}$, the k most likely next-tokens under M_{new} .

Because the models operate over different vocabularies, we define a decode–retokenize mapping $f : \mathcal{V}_{\text{new}} \rightarrow \mathcal{V}_{\text{old}}$ that projects each candidate token from the new model into the clinical model’s vocabulary by decoding a token from \mathcal{V}_{new} to its string form and re-tokenizing the string using \mathcal{T}_{old} . For $i \in \mathcal{C}_t$, we define $f(i)$ as the first non-space token in $\mathcal{T}_{\text{old}}(\text{decode}_{\text{new}}(i))$. For each candidate token $i \in \mathcal{C}_t$, we update its log-probability:

$$s(i) = \log p_{\text{new}}(i | \mathbf{x}_{1:t}) + \alpha \left(\log p_{\text{old-clin}}(f(i) | \mathbf{x}_{1:t}) - \log p_{\text{old}}(f(i) | \mathbf{x}_{1:t}) \right)$$

which adds a correction representing the isolated clinical-domain signal. The selected token, $i^* = \arg \max_{i \in \mathcal{C}_t} s(i)$, is appended to the context. We set $k = 20$, reflecting the top-20 log-probabilities exposed by black-box API providers (OpenAI, 2025; Google DeepMind, 2025) and fix $\alpha = 1.0$. See Appendix B for hyperparameter analysis.

3.2. Models and Baselines

We use Qwen3-30B (Yang et al., 2025) as the new-generation general-domain model (M_{new}) and MeLLaMA-13B-chat (Xie et al., 2024) as the old-generation clinical model ($M_{\text{old-clin}}$). MeLLaMA-13B-chat is continually pre-trained from LLaMA-2-13B-base (M_{old}) and instruction-tuned on clinical notes and tasks. We compare CAPT to proxy tuning (Liu et al., 2024) and UNiTE (Yao et al., 2025), the strongest probability-level ensembling method to date. Since proxy tuning requires a shared tokenizer between models, we use LLaMA-2-70B-chat ($M_{\text{old-L}}$) (Touvron et al., 2023) as its general-domain model.

3.3. Evaluation Tasks

As shown in Table 1, we evaluate CAPT on six clinical classification and text-generation tasks to assess

Figure 2: **CAPT decoding process.** At each step, CAPT (1) selects top- k token candidates from the new-generation general domain model, (2) maps each token candidate to the clinical model’s vocabulary via retokenization, (3) computes the clinical knowledge offset as the log-probability difference between the clinical model and its untuned counterpart, and (4) re-ranks candidates by adding the clinical offset to the original log-probabilities. In this example, CAPT selects “perf” (perfusion) over “cardiac” for monitoring a forearm graft, as the clinical offset adjusts the ranking to favor the contextually appropriate term. Full example context can be found in Figure 4a.

using the clinical model’s tokenizer to obtain a corresponding clinical token. We then adjust the candidate’s log-probability by adding the log-probability difference between the clinical model and its base counterpart. The highest-scoring token after adjustment is appended to the context.

Table 1: Evaluation Tasks

Task	Input	Output
Classification		
MedNLI	Premise from clinical note + hypothesis	Entailment/Neutral/Contradiction
MTS-Specialty	Medical transcription note	One of 40 medical specialties
CLIP	Discharge summary sentence	One or more of 7 follow-up labels
Text Generation		
MIMIC-RRS	Radiology report findings section	Impression section
MIMIC-BHC	Full discharge summary note	Brief Hospital Course summary
MTS-Procedure	Operative note or surgical case description	Postoperative treatment plan

performance across a broad range of use cases. For classification tasks, we report Macro-F1 and accuracy on 200 random test samples from: *MedNLI* (Romanov and Shivade, 2018), *MTS-Specialty* (MTSsamples, 2025), and *CLIP* (Mullenbach et al., 2021).

For text generation, we report MedHELM LLM-jury scores (Bedi et al., 2025) and MedVAL percent of risk-free outputs (Aali et al., 2025) on 100 samples from: *MIMIC-BHC* (Aali et al., 2024), *MTS-Procedures* (MTSsamples, 2025), and *MIMIC-RRS* (Chen et al., 2023).

MedHELM’s LLM-jury metric evaluates accuracy, completeness, and clarity of generated outputs against reference outputs. Our results can be compared to MedHELM’s leaderboard, and CAPT can be readily applied to any of the 37 MedHELM clinical classification, text generation, or question-answering tasks. MedVAL is a fine-tuned LLM that assigns a risk score to outputs, where level 1 is risk-free and level 4 is the highest risk level. Both metrics are validated against clinician judgement. See Appendix A for details on experimental setup and task prompts.

3.4. Token-level Analysis

To understand how CAPT leads to token-level changes, we manually categorized 280 of the top 25% most frequently generated tokens in the MTS-Procedure task into semantic categories. The remaining tokens were excluded because they were ambiguous, subword fragments, or otherwise not semantically interpretable. We use MTS-Procedure because it is publicly available; our other text-generation tasks rely on MIMIC (Johnson et al., 2016), whose raw data cannot be shared. See Appendix C for token categories and corresponding tokens.

For each category, we report the mean per-token log-probability difference between the clinical model

$M_{old-clin}$ and its base counterpart M_{old} . Because shifts are computed over generated tokens, positive values indicate tokens preferentially up-weighted by the clinical model, while negative values indicate tokens supported by the new-generation general-domain model, such that they remain selected even when the clinical offset disfavors them.

3.5. Physician Case Study

We present five randomly selected full-length CAPT outputs from the MTS-Procedure task, annotated with token-level shifts and expert qualitative analysis. The analysis was conducted by two board-certified physicians, an interventional radiologist with 30 years of clinical experience and an anesthesiologist with 5 years of clinical experience and an expertise in clinical language models. Similar to the token-level analysis, we use MTS-Procedure because it is publicly available, whereas our other text-generation tasks rely on MIMIC (Johnson et al., 2016). The physicians assessed clinical accuracy, appropriateness, and overall utility.

4. Results

4.1. Model Performance

CAPT outperforms baselines and existing model-ensembling approaches. Table 2 summarizes performance across clinical classification and text generation tasks. M_{new} outperforms $M_{old-clin}$ by an average of 53.6% in Macro-F1 on classification benchmarks and 52.6% in LLM-jury scores on text-generation benchmarks. This performance gap likely explains why UNiTE, which averages the probability scores of $M_{old-clin}$ and M_{new} , performs worse than M_{new} alone by 8-28% on average. Proxy tuning improves

Table 2: **Performance comparison across tasks and methods.** Bold indicates best performance. F1 = Macro-F1, Acc = Accuracy, LLM-J = MedHelm LLM-jury score (average out of 5), %RF = MedVAL % of risk-free outputs.

Method / Models	Classification						Text-Generation					
	MedNLI		MTS-Specialty		CLIP		MTS-Procedure		MIMIC-RRS		MIMIC-BHC	
	F1	Acc	F1	Acc	F1	Acc	LLM-J	%RF	LLM-J	%RF	LLM-J	%RF
Old-gen large general (M_{old-L}) LLaMA-2-70B-chat	0.598	0.635	0.094	0.280	0.097	0.205	3.214	26.67	2.453	57.50	3.323	13.00
Old-gen clinical ($M_{old-clin}$) MeLLaMA-13B-chat	0.539	0.535	0.080	0.210	0.192	0.785	1.989	45.83	3.379	35.00	2.906	14.14
New-gen general (M_{new}) Qwen3-30B	0.858	0.850	0.140	0.350	0.240	0.540	3.825	63.33	4.394	65.00	3.932	30.00
Proxy Tuning $M_{old-L} + M_{old-clin} + M_{old}$	0.575	0.620	0.095	0.275	0.143	0.420	3.466	55.00	3.898	34.17	3.456	8.00
UniTE $M_{new} + M_{old-clin}$	0.873	0.875	0.060	0.100	0.217	0.560	3.486	50.00	3.903	46.67	3.747	20.00
CAPT $M_{new} + M_{old-clin} + M_{old}$	0.886	0.885	0.135	0.320	0.224	0.580	3.882	70.83	4.414	56.67	3.973	31.00

upon both M_{old-L} and $M_{old-clin}$ despite their large performance gap, demonstrating the potential of contrastive decoding-based methods to bridge models of distinct capabilities. However, M_{new} consistently outperforms proxy tuning, highlighting how the method’s reliance on shared vocabularies limits its applicability as new base models are released.

CAPT outperforms M_{new} on 8 of 12 metrics, while UniTE improves on only 2. Averaged across all tasks and metrics, CAPT outperforms UniTE and proxy tuning by 17.6% and 41.4%, respectively. CAPT’s gains are particularly pronounced on text-generation tasks, where clinical reasoning and terminology precision are critical. For these tasks, CAPT yields substantially more risk-free outputs, improving by 13.94 and 20.44 percentage points over UniTE and proxy tuning respectively. This indicates meaningful improvements in clinical safety and utility alongside performance gains.

4.2. Token-level Analysis

CAPT selectively integrates clinical domain knowledge. Figure 3 summarizes mean log-probability shifts by semantic token category. Positive values indicate categories for which CAPT increases token preference via the clinical-model offset between $M_{old-clin}$ and M_{old} , while negative values indicate categories for which the offset decreases preference and M_{new} retains stronger influence.

$M_{old-clin}$ strongly influences tokens associated with clinical decision-making (e.g., *Clinical Decision Action Headers*, *Clinical Assessment Terms*) and documentation style (e.g., *Clinical Reporting Verbs*, *Clin-*

ical Hedging, *Condition State Descriptors*). These positive shifts reflect the clinical model’s learned conventions for structuring clinical assessments, expressing uncertainty, and framing care decisions. These stylistic patterns are largely absent from general-domain pretraining corpora, which typically exclude clinical notes, indicating that CAPT effectively transfers knowledge learned by an older-generation model.

In contrast, M_{new} dominates linguistic structure and formatting (e.g., *Formatting Tokens*, *General Morphemes*, *Medical Suffixes*), reflecting that the general-domain model controls grammatical coherence and document structure. Similarly, *Time Demographic Modifiers* show negative shifts, indicating that M_{new} influences temporal precision and reasoning. These patterns demonstrate that CAPT preserves the newer model’s capabilities in general reasoning and fluency.

Notably, medical knowledge categories exhibit near-zero shifts, indicating both models converge on similar preferences for foundational medical concepts. Core medical terminology including *Gynecologic Terms*, *Diagnoses*, *Medical Roots*, and *Physiological Progress* all cluster near zero. This suggests CAPT primarily affects how M_{new} ’s medical knowledge is expressed to match clinical documentation conventions.

4.3. Physician Case Study

Figure 4 shows samples from two case studies generating postoperative treatment plans for the MTSamples-Procedure task. See Appendix D for ex-

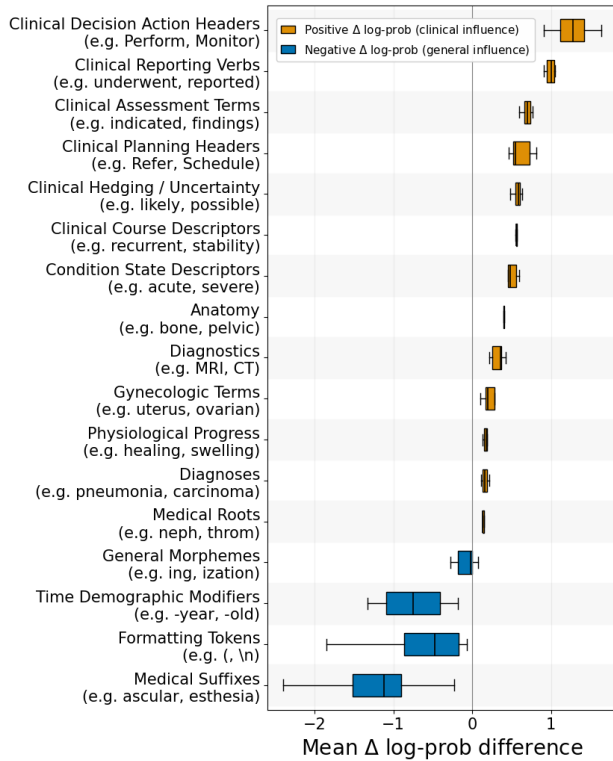


Figure 3: **Mean log-probability offset between $M_{old-clin}$ and M_{old} of generated tokens by semantic category.** Positive shifts (orange) indicate increased influence of $M_{old-clin}$, while negative shifts (blue) of M_{new} . All token categories are shown in Appendix C, Figure 5.

tended analysis, LLM jury evaluation, and additional examples.

4.3.1. CASE STUDY 1

Figure 4a shows a sample of a CAPT-generated treatment plan for a forearm arteriovenous graft procedure (graft connecting vein to artery).

CAPT token-level shifts improve clinical accuracy. CAPT replaces “72” with “48” in the postoperative monitoring timeline, narrowing the range from 24-72 hours to 24-48 hours in line with typical discharge practices for this procedure. Similarly, M_{new} recommends monitoring “cardiac” symptoms, which are not relevant to a forearm graft; CAPT correctly emphasizes perfusion (i.e. blood flow).

CAPT also shifts generation toward more precise clinical terminology, resulting in a more clinically actionable and descriptive plan. CAPT replaces M_{new} ’s “neuro” token, which likely implies “neuro check,” with “dist.” This change directs generation towards monitoring “distal limb circulation,” which encompasses relevant components of the neurologic assessment (e.g., sensation and motor function) as well as other important checks (e.g., capillary refill). Similarly, CAPT inserts “AV” before “graft” to clarify the graft type and replaces “agents” with “analgesics” to specify pain medication rather than any generic medication class. CAPT also heavily preferences “proph” for “prophylactic antibiotics” after “Administer,” specifying the preventative nature of the antibiotic administration.

One limitation is that CAPT preserved over-explanatory content intended for non-expert audiences (e.g., naming specific antibiotics), which are unnecessary for nurses responsible for post-operative management. These additions often appear in parentheses (e.g., “(”) that are favored by M_{new} , reflecting its general-audience training objective.

4.3.2. CASE STUDY 2

Figure 4b shows a CAPT-generated treatment plan for a low transverse cesarean section (C-section), a surgical procedure in which the fetus is delivered through a horizontal lower uterine incision.

CAPT increases procedure-specific, safety-critical recommendations. The clinical model increases preference towards “erythema, discharge, and dehiscence,” procedure-specific signs of infection critical to monitor in C-sections. CAPT also adds “fundal” height, a highly procedure-specific physical exam sign that serves as a proxy for adequate uterine tone (contraction of the uterus, which is essential for hemorrhage prevention). CAPT replaces “foul” with “sudden” when describing changes in lochia (normal bloody discharge following delivery), reflecting a shift from monitoring general infection signs to detecting postpartum hemorrhage. Sudden increases in lochia are a sign of delayed postpartum hemorrhage, whereas foul odor is a sign of infection or chorioamnionitis, a more rare complication after delivery than postpartum hemorrhage. Similarly, CAPT replaces monitoring for “fever” with “passage of large clots,” another sign of postpartum hemorrhage. Overall, this shift in focus towards postpartum hemorrhage is clinically appropriate, as hemorrhage is the leading cause

a) Prompt – Forearm arteriovenous graft

Here are information about a patient, return a reasonable treatment plan for the patient. Patient Notes: Medical Specialty: Nephrology. Sample Name: Fogarty thrombectomy. Description: Fogarty thrombectomy, left forearm arteriovenous Gore-Tex bridge fistula and revision of distal anastomosis with 7 mm interposition Gore-Tex graft. Chronic renal failure and thrombosed left forearm arteriovenous Gore-Tex bridge fistula. PREOPERATIVE/POSTOPERATIVE DIAGNOSES: 1. Chronic renal failure. 2. Thrombosed left forearm arteriovenous Gore-Tex bridge fistula. PROCEDURE PERFORMED: 1. Fogarty thrombectomy, left forearm arteriovenous Gore-Tex bridge fistula. 2. Revision of distal anastomosis with 7 mm interposition Gore-Tex graft. ANESTHESIA: General with controlled ventilation.

CAPT-generated output

Treatment Plan: **Post-Operative Management Following Fogarty Thrombectomy and AV Graft Revision**
1. Immediate Post-Operative Care (First 24-48 Hours):
 - **Monitor** vital signs (**cardiac perfusion** and **neuro distal limb circulation (capillary refill, temperature, color, pulses sensation, pulse motor function)**).
 - Monitor for **laboratory** values: CBC, electrolytes, **renal function (BUN, creatinine)**, coagulation profile (PT/INR, aPTT).
 - **Ensure** **Administer** **prophylactic** antibiotics (if indicated (e.g., cefazolin or equivalent) **per institutional surgical protocol** to prevent **graft infection**, especially typically for 24-48 hours **post-op** unless contraindicated).
 - **Provide** **Pain control management** with non-nephrotoxic **agents (analgesics)** (e.g., acetaminophen); **avoid NSAIDs** due to risk of renal impairment **and**.
■ Positive Δ log-prob (clinical influence) ■ Negative Δ log-prob (general influence)

b) Prompt – Low transverse C-section

Here are information about a patient, return a reasonable treatment plan for the patient. Patient Notes: Medical Specialty: Obstetrics / Gynecology Sample Name: Low transverse C-Section Description: Primary low-transverse cesarean section. PREOPERATIVE/POSTOPERATIVE DIAGNOSES: 1. Intrauterine pregnancy at term. 2. Arrest of dilation. PROCEDURE PERFORMED: Primary low-transverse cesarean section. ANESTHESIA:

CAPT-generated output

4. Wound Care and Infection Prevention:
 - ~~Ins~~ **Keep** surgical incision clean and **dry**; **inspect** daily for signs of infection (**erythema, swelling, discharge, dehiscence**).
 - ~~Prophylactic~~ **antibiotics** ~~were~~ likely administered intraoperatively (no further **additional antibiotics prophylaxis** needed unless **indicated** by clinical **signs findings**).
 - **Encourage** early ambulation to **reduce promote healing** circulation and **reduce** risk of **infection thromboembolism**.
5. Uterine Involution and Lochia Monitoring Management:
 - **Monitor** **fundal** height and **ut consistency** (uterus **ensure** the uterus is firm and midline **to prevent** postpartum hemorrhage).
 - **Assess** lochia (**color amount, color, odor**) daily; **report** any **foot sudden** increase in flow, foul odor, or **fever passage of large clots**.
 - **Encourage** breastfeeding **to or milk pumping** to stimulate uterine **contractions** and aid involution.

Figure 4: CAPT-generated post-operative management plans for two MTS-Procedure examples. Orange and blue highlights indicate stronger influence from $M_{old-clin}$ and M_{new} , respectively. Bolded tokens denote when the top-choice token changed after CAPT adjustment, with the original top-choice crossed out. Green circles mark tokens discussed in main text. The full outputs, extended analysis, LLM jury evaluation, and additional case studies are shown in Appendix D.

of maternal morbidity worldwide and undergoing a C-section increases a patient’s risk (Yunas et al., 2025).

CAPT additionally reduces clinical errors. In the recommendation to encourage early ambulation (moving during recovery for circulation), CAPT replaces “infection” with “thromboembolism” (blood clot forming and traveling to block blood flow). Early ambulation reduces risk of thromboembolism, but does not reduce risk of infection.

4.3.3. TAKEAWAYS

Across both cases, CAPT makes three types of improvements: (1) *increased clinical specificity* (replacing generic term “agents” with more precise “analgesics”), (2) *enhanced clinical and contextual accuracy* (“24-72” → “24-48” for postoperative monitoring timeline; “infection” → “thromboembolism” for

risks reduced with ambulation), and (3) *emphasizing procedure-specific, safety-critical monitoring* (“distal limb circulation” for a lower forearm graft; postpartum hemorrhage signs for a C-section). We find that CAPT preserves M_{new} ’s explanatory style, which may be redundant for expert clinicians. Although a limitation, this suggests that CAPT incorporates clinical expertise without altering the base model’s communication patterns.

5. Discussion

We introduce *Cross-Architecture Proxy Tuning* (CAPT), a training-free model-ensembling strategy that adapts newly released general-domain language models to clinical use by reusing legacy clinical models. CAPT targets a practical bottleneck in clinical

NLP deployment: continued pretraining and instruction tuning must be repeated for each new model generation, which is often infeasible for resource-constrained health systems. By operating at decoding time and supporting disjoint tokenizers, CAPT provides a path to benefit from new architecture advances without reinvesting in large-scale clinical re-training.

Across six clinical classification and text-generation tasks, CAPT consistently improves over strong baselines and prior heterogeneous ensembling methods, including cases where naive averaging harms performance due to large capability gaps between models. Our mechanistic analyses offer evidence for how CAPT achieves these gains. Token-level shifts indicate that the clinical model primarily influences decision-relevant phrasing and documentation conventions, while the new-generation model retains control over linguistic structure and global coherence. This division of labor is consistent with the intended design of CAPT: rather than transferring broad medical knowledge, CAPT preferentially alters how medical knowledge is expressed so that outputs better match clinical note style and clinical action framing. Physician case studies support this interpretation, showing that CAPT promotes reduces context errors that can affect downstream usability.

This study has several limitations that motivate future work. First, clinically meaningful improvements can be sparse and context-dependent, and therefore are not always well captured by quantitative metrics. While we partially address this gap via physician review of CAPT outputs, robust evaluation of opened clinical generation remains an open challenge, particularly for measuring rare but safety-critical failures and subtle improvements in contextual appropriateness. Second, CAPT can inherit communication patterns from the general-domain model, including over-explanatory content that may be unnecessary for expert clinical audiences. Adaptively scaling the clinical-model adjustment coefficient (i.e. α) based on intended user expertise could address this limitation. Third, CAPT increases inference-time compute by requiring forward passes through three models. Although we did not observe a significant wall-clock latency difference between Qwen3-30b alone and CAPT (which additionally includes MeLLaMA-13b-chat and LLaMA-2-13-base), the additional memory requirements may matter in deployment settings with stricter constraints. In this study,

we use MeLLaMA-13B models, which are the smallest models in the MeLLaMA family and are state-of-the-art among open-source clinical language models. Exploring CAPT with smaller or more specialized clinical models may further improve the feasibility of this approach.

Together, these results position CAPT as a practical bridge between rapidly advancing general-domain model architectures and slower-moving clinical adaptation pipelines. CAPT enables health systems and clinical NLP teams to reuse prior investments in domain-adapted models to make new-generation models more clinically aligned without additional training. More broadly, this work suggests a deployment strategy where legacy clinical models serve as reusable, modular adapters that can be paired with future foundation models, accelerating translation of new capabilities into clinical workflows while keeping adaptation costs and governance burdens manageable.

References

- Asad Aali, Dave Van Veen, Yamin Ishraq Arefeen, Jason Hom, Christian Bluethgen, Eduardo Pontes Reis, Sergios Gatidis, Namuun Clifford, Joseph Daws, Arash S Tehrani, Jangwon Kim, and Akshay S Chaudhari. A dataset and benchmark for hospital course summarization with adapted large language models. *Journal of the American Medical Informatics Association*, 32(3):470–479, December 2024. ISSN 1527-974X. doi: 10.1093/jamia/ocae312. URL <http://dx.doi.org/10.1093/jamia/ocae312>.
- Asad Aali, Vasiliki Bikia, Maya Varma, Nicole Chiou, Sophie Ostmeier, Arnav Singhvi, Magdalini Paschali, Ashwin Kumar, Andrew Johnston, Karimar Amador-Martinez, Eduardo Juan Perez Guerrero, Paola Naovi Cruz Rivera, Sergios Gatidis, Christian Bluethgen, Eduardo Pontes Reis, Eddy D. Zandee van Rilland, Poonam Laxmappa Hosamani, Kevin R Keet, Minjoung Go, Evelyn Ling, David B. Larson, Curtis Langlotz, Roxana Daneshjou, Jason Hom, Sanmi Koyejo, Emily Alsentzer, and Akshay S. Chaudhari. Medval: Toward expert-level medical text validation with language models, 2025. URL <https://arxiv.org/abs/2507.03152>.
- David A. Alber, Zeming Yang, Alexey Alyakin, et al. Medical large language models are

- vulnerable to data-poisoning attacks. *Nature Medicine*, 31:618–626, 2025. doi: 10.1038/s41591-024-03445-1. URL <https://doi.org/10.1038/s41591-024-03445-1>.
- Emily Alsentzer, John R. Murphy, Willie Boag, Wei-Hung Weng, Di Jin, Tristan Naumann, and Matthew B. A. McDermott. Publicly available clinical bert embeddings, 2019. URL <https://arxiv.org/abs/1904.03323>.
- Elham Asgari, Nina Montaña-Brown, Magda Dubois, Saleh Khalil, Jasmine Balloch, Joshua Au Yeung, and Dominic Pimenta. A framework to assess clinical safety and hallucination rates of llms for medical text summarisation. *npj Digital Medicine*, 8(1): 274, May 2025. URL <https://www.nature.com/articles/s41746-025-01670-7>.
- Suhana Bedi, Hejie Cui, Miguel Fuentes, Alyssa Unell, Michael Wornow, Juan M. Banda, Nikesh Kotecha, Timothy Keyes, Yifan Mai, Mert Oez, Hao Qiu, Shrey Jain, Leonardo Schettini, Mehr Kashyap, Jason Alan Fries, Akshay Swaminathan, Philip Chung, Fateme Nateghi, Asad Aali, Ashwin Nayak, Shivam Vedak, Sneha S. Jain, Birju Patel, Oluseyi Fayanju, Shreya Shah, Ethan Goh, Dong han Yao, Brian Soetikno, Eduardo Reis, Sergio Gatidis, Vasu Divi, Robson Capasso, Rachna Saralkar, Chia-Chun Chiang, Jenelle Jindal, Tho Pham, Faraz Ghoddsusi, Steven Lin, Albert S. Chiou, Christy Hong, Mohana Roy, Michael F. Gensheimer, Hinesh Patel, Kevin Schulman, Dev Dash, Danton Char, Lance Downing, Francois Grolleau, Kameron Black, Bethel Mieso, Aydin Zahedivash, Wen wai Yim, Harshita Sharma, Tony Lee, Hannah Kirsch, Jennifer Lee, Nerissa Ambers, Carlene Lugtu, Aditya Sharma, Bilal Mawji, Alex Alekseyev, Vicky Zhou, Vikas Kakkar, Jarrod Helzer, Anurang Revri, Yair Barnett, Roxana Daneshjou, Jonathan Chen, Emily Alsentzer, Keith Morse, Nirmal Ravi, Nima Aghaeepour, Vanessa Kennedy, Akshay Chaudhari, Thomas Wang, Sanmi Koyejo, Matthew P. Lungren, Eric Horvitz, Percy Liang, Mike Pfeffer, and Nigam H. Shah. Medhelm: Holistic evaluation of large language models for medical tasks, 2025. URL <https://arxiv.org/abs/2505.23802>.
- Zhihong Chen, Maya Varma, Xiang Wan, Curtis Langlotz, and Jean-Benoit Delbrouck. Toward expanding the scope of radiology report summarization to multiple anatomies and modalities. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, page 469–484. Association for Computational Linguistics, 2023. doi: 10.18653/v1/2023.acl-short.41. URL <http://dx.doi.org/10.18653/v1/2023.acl-short.41>.
- Google DeepMind. Gemini 3 pro model card, November 2025. URL <https://deepmind.google/models/model-cards/gemini-3-pro/>. Accessed: 2025-01.
- Paul Hager, Friederike Jungmann, Robbie Holland, Kunal Bhagat, Inga Hubrecht, Manuel Knauer, Jakob Vielhauer, Marcus Makowski, Rickmer Braren, Georgios Kaissis, and Daniel Rueckert. Evaluation and mitigation of the limitations of large language models in clinical decision-making. *Nature Medicine*, 30:2613–2622, July 2024. doi: 10.1038/s41591-024-03068-6. URL <https://www.nature.com/articles/s41591-024-03068-6>.
- Tianyu Han, Lisa C. Adams, Jens-Michalis Papaioannou, Paul Grundmann, Tom Oberhauser, Alexei Figueroa, Alexander Löser, Daniel Truhn, and Keno K. Bressen. Medalpaca – an open-source collection of medical conversational ai models and training data, 2025. URL <https://arxiv.org/abs/2304.08247>.
- Yichong Huang, Xiaocheng Feng, Baohang Li, Yang Xiang, Hui Wang, Bing Qin, and Ting Liu. Ensemble learning for heterogeneous large language models with deep parallel collaboration, 2024. URL <https://arxiv.org/abs/2404.12715>.
- L. Y. Jiang, X. C. Liu, N. P. Nejatian, et al. Health system-scale language models are all-purpose prediction engines. *Nature*, 619:357–362, 2023. doi: 10.1038/s41586-023-06160-y. URL <https://doi.org/10.1038/s41586-023-06160-y>.
- Alistair E. W. Johnson, Tom J. Pollard, Lu Shen, Liwei H. Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G. Mark. Mimic-iii, a freely accessible critical care database. *Scientific Data*, 3:160035, 2016. doi: 10.1038/sdata.2016.35. URL <https://doi.org/10.1038/sdata.2016.35>.
- Jonathan Kim, Anna Podlasek, Kie Shidara, Feng Liu, Ahmed Alaa, and Danilo Bernardo. Limitations of large language models in clinical problem-

- solving arising from inflexible reasoning. *Scientific Reports*, 15(1), November 2025. ISSN 2045-2322. doi: 10.1038/s41598-025-22940-0. URL <http://dx.doi.org/10.1038/s41598-025-22940-0>.
- James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A. Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, Demis Hassabis, Claudia Clopath, Dharshan Kumaran, and Raia Hadsell. Overcoming catastrophic forgetting in neural networks. *Proceedings of the National Academy of Sciences*, 114(13):3521–3526, March 2017. ISSN 1091-6490. doi: 10.1073/pnas.1611835114. URL <http://dx.doi.org/10.1073/pnas.1611835114>.
- Eric Lehman and Alistair Johnson. Clinical-t5: Large language models built using mimic clinical text. PhysioNet (version 1.0.0), 2023. URL <https://doi.org/10.13026/rj8x-v335>. RRID:SCR_007345.
- Eric Lehman, Evan Hernandez, Diwakar Mahajan, Jonas Wulff, Micah J. Smith, Zachary Ziegler, Daniel Nadler, Peter Szolovits, Alistair Johnson, and Emily Alsentzer. Do we still need clinical language models?, 2023. URL <https://arxiv.org/abs/2302.08091>.
- Xiang Lisa Li, Ari Holtzman, Daniel Fried, Percy Liang, Jason Eisner, Tatsunori Hashimoto, Luke Zettlemoyer, and Mike Lewis. Contrastive decoding: Open-ended text generation as optimization. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12286–12312, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.687. URL <https://aclanthology.org/2023.acl-long.687/>.
- A. Liu, X. Han, Y. Wang, Y. Tsvetkov, Y. Choi, and N. A. Smith. Tuning language models by proxy. *arXiv preprint arXiv:2401.08565*, 2024.
- Alisa Liu, Maarten Sap, Ximing Lu, Swabha Swayamdipta, Chandra Bhagavatula, Noah A. Smith, and Yejin Choi. DExperts: Decoding-time controlled text generation with experts and anti-experts. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli, editors, *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 6691–6706, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-long.522. URL <https://aclanthology.org/2021.acl-long.522/>.
- Eric Mitchell, Rafael Rafailov, Archit Sharma, Chelsea Finn, and Christopher D. Manning. An emulator for fine-tuning large language models using small language models. *arXiv preprint arXiv:2310.12962*, 2023.
- MTSamples. Mtsamples: Medical transcription samples. <https://www.mtsamples.com/>, 2025. Accessed 2025-09-08.
- James Mullenbach, Yada Pruksachatkun, Sean Adler, Jennifer Seale, Jordan Swartz, T. Greg McKelvey, Hui Dai, Yi Yang, and David Sontag. Clip: A dataset for extracting action items for physicians from hospital discharge notes, 2021. URL <https://arxiv.org/abs/2106.02524>.
- OpenAI. GPT-5 System Card. Technical report, OpenAI, August 2025. URL <https://cdn.openai.com/gpt-5-system-card.pdf>. Accessed: 2025-11-03.
- Alexey Romanov and Chaitanya Shivade. Lessons from natural language inference in the clinical domain, 2018. URL <https://arxiv.org/abs/1808.06752>.
- Karan Singhal, Shekoofeh Azizi, Tu Tu, et al. Large language models encode clinical knowledge. *Nature*, 620:172–180, 2023a. doi: 10.1038/s41586-023-06291-2.
- Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Le Hou, Kevin Clark, Stephen Pföhl, Heather Cole-Lewis, Darlene Neal, Mike Schaeckermann, Amy Wang, Mohamed Amin, Sami Lachgar, Philip Mansfield, Sushant Prakash, Bradley Green, Ewa Dominowska, Blaise Agueray Arcas, Nenad Tomasev, Yun Liu, Renee Wong, Christopher Semturs, S. Sara Mahdavi, Joelle Barral, Dale Webster, Greg S. Corrado, Yossi Matias, Shekoofeh Azizi, Alan Karthikesalingam, and Vivek Natarajan. Towards expert-level medical question answering with large language mod-

- els, 2023b. URL <https://arxiv.org/abs/2305.09617>.
- Vinith M. Suriyakumar, Ayush Sekhari, and Ashia Wilson. Ucd: Unlearning in llms via contrastive decoding, 2025. URL <https://arxiv.org/abs/2506.12097>.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models, 2023. URL <https://arxiv.org/abs/2307.09288>.
- Hanyin Wang, Chufan Gao, Bolun Liu, Qiping Xu, Guleid Hussein, Mohamad El Labban, Kingsley Iheasirim, Hariprasad Korsapati, Chuck Outcalt, and Jimeng Sun. Towards adapting open-source large language models for expert-level clinical note generation, 2025. URL <https://arxiv.org/abs/2405.00715>.
- Eric Wu, Kevin Wu, and James Zou. Limitations of learning new and updated medical knowledge with commercial fine-tuning large language models. *NEJM AI*, 2(8):AIcs2401155, 2025a. doi: 10.1056/AIcs2401155. URL <https://ai.nejm.org/doi/full/10.1056/AIcs2401155>.
- Jiageng Wu, Bowen Gu, Ren Zhou, Kevin Xie, Doug Snyder, Yixing Jiang, Valentina Carducci, Richard Wyss, Rishi J Desai, Emily Alsentzer, Leo Anthony Celi, Adam Rodman, Sebastian Schneeweiss, Jonathan H. Chen, Santiago Romero-Brufau, Kueiyu Joshua Lin, and Jie Yang. Bridge: Benchmarking large language models for understanding real-world clinical practice text, 2025b. URL <https://arxiv.org/abs/2504.19467>.
- Qianqian Xie, Qingyu Chen, Aokun Chen, Cheng Peng, Yan Hu, Fongci Lin, Xueqing Peng, Jimin Huang, Jeffrey Zhang, Vipina Keloth, Xinyu Zhou, Lingfei Qian, Huan He, Dennis Shung, Lucila Ohno-Machado, Yonghui Wu, Hua Xu, and Jiang Bian. Me llama: Foundation large language models for medical applications, 2024. URL <https://arxiv.org/abs/2402.12749>.
- Yangyifan Xu, Jinliang Lu, and Jiajun Zhang. Bridging the gap between different vocabularies for llm ensemble, 2024. URL <https://arxiv.org/abs/2404.09492>.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujiu Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiayi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. Qwen3 technical report, 2025. URL <https://arxiv.org/abs/2505.09388>.
- Xi Yang, Aokun Chen, Nima PourNejatian, Hoo Chang Shin, Kaleb E Smith, Christopher Parisien, Colin Compas, Cheryl Martin, Mona G Flores, Ying Zhang, Tanja Magoc, Christopher A Harle, Gloria Lipori, Duane A Mitchell, William R Hogan, Elizabeth A Shenkman, Jiang Bian, and Yonghui Wu. Gatortron: A large clinical language model to unlock patient information from unstructured electronic health records, 2022. URL <https://arxiv.org/abs/2203.03540>.
- Yuxuan Yao, Han Wu, Mingyang Liu, Sichun Luo, Xiongwei Han, Jie Liu, Zhijiang Guo, and Linqi

Song. Determine-then-ensemble: Necessity of top-k union for large language model ensembling, 2025. URL <https://arxiv.org/abs/2410.03777>.

Wen Yim, Yujuan Fu, Asma Ben Abacha, Neal Snider, Thomas Lin, and Meliha Yetisgen. Acibench: a novel ambient clinical intelligence dataset for benchmarking automatic visit note generation, 2023. URL <https://arxiv.org/abs/2306.02022>.

Idnan Yunas, Md Asiful Islam, Kulandaipalayam N. Sindhu, Adam J. Devall, Marcelina Podeseck, Sayeda Sadia Alam, Shoumik Kundu, Kristie-Marie Mammoliti, Ashraf Aswat, Malcolm J. Price, Javier Zamora, Olufemi T. Oladapo, Ioannis Gallos, and Arri Coomarasamy. Causes of and risk factors for postpartum haemorrhage: a systematic review and meta-analysis. *The Lancet*, 405(10488): 1468–1480, 2025. doi: 10.1016/S0140-6736(25)00448-9. Epub ahead of print: Apr 3, 2025.

T. Zack, E. Lehman, M. Suzgun, J. A. Rodriguez, L. A. Celi, J. Gichoya, D. Jurafsky, P. Szolovits, D. W. Bates, R. E. Abdunour, A. J. Butte, and E. Alsentzer. Assessing the potential of GPT-4 to perpetuate racial and gender biases in health care: A model evaluation study. *The Lancet Digital Health*, 6(1):e12–e22, January 2024. doi: 10.1016/S2589-7500(23)00225-X.

Hongbo Zhang, Junying Chen, Feng Jiang, Fei Yu, Zhihong Chen, Jianquan Li, Guiming Chen, Xiangbo Wu, Zhiyi Zhang, Qingying Xiao, Xiang Wan, Benyou Wang, and Haizhou Li. Huatuogpt, towards taming language model to be a doctor, 2023. URL <https://arxiv.org/abs/2305.15075>.

Xinlu Zhang, Chenxin Tian, Xianjun Yang, Lichang Chen, Zekun Li, and Linda Ruth Petzold. Alpacare:instruction-tuned large language models for medical application, 2025. URL <https://arxiv.org/abs/2310.14558>.

Appendix A. Experimental Setup

For classification tasks, we use prefix-constrained generation to constrain outputs to an explicit JSON schema consisting of a “reason” field (a free-text string with a maximum length of 600 characters) and a “label” field (a string or array of strings drawn from a predefined, task-specific label set). At each decoding step, the current prefix is used to determine the

set of tokens that can legally follow while still permitting completion of a valid JSON object; all other tokens are masked out. In CAPT, the constraint is applied to the base model’s token distribution, since the base model alone determines the generated tokens. In proxy tuning and UNiTE, the constraint is applied to the token candidates considered during model combination, ensuring that only schema-valid tokens remain eligible after score aggregation. We employ constrained decoding to control for performance differences attributable to schema-adherence during instruction following.

Table 6 contains a description of each task and example prompt. All models were 4-bit quantized to fit on NVIDIA A100 40GB GPUs, such that our CAPT setup could run on a single 40GB GPU.

Appendix B. Hyperparameter Analysis

To evaluate the effect of k and α , we conducted experiments on two held-out, validation datasets. We report Macro-F1 on the validation split of the *MedNLI* classification task, and MedHelm LLM-jury score on *ACI-Bench*, a text-generation task of generating a structured clinical note from patient-doctor dialogues (Yim et al., 2023). We use *ACI-Bench* because MedHelm does not provide a validation set for its text-generation tasks and we did not want to reduce sample size of our tasks. Table 3 shows that varying k has minimal impact on performance. We select $k = 20$ because black-box API providers provide the top-20 logprobs for each token generation (OpenAI, 2025; Google DeepMind, 2025). We set $\alpha = 1.0$ because it achieves best performance on both tasks, as shown by Table 4.

k	MedNLI (Macro-F1)	ACI-Bench (LLM-jury score)
5	0.8944	4.377
10	0.8944	4.380
15	0.8944	4.377
20	0.8992	4.379

Table 3: Top- k parameter experiments.

α	MedNLI (Macro-F1)	ACI-Bench (LLM-jury score)
0.5	0.8486	4.360
0.7	0.8537	4.379
1.0	0.8619	4.621

Table 4: α parameter analysis.

Appendix C. Token-Level Analysis

Tables 7 and 8 show token categories and corresponding tokens. Figure 5 shows the full version of Figure 3, with all token categories shown. Due space constraints, Figure 3 only includes the categories discussed in the main text.

Appendix D. Physician Case Study

Figures 6, 9, 12, 15, and 18 contain the annotated CAPT outputs. For each case, we provide the LLM jury evaluation (Figures 7, 10, 13, 16, 19) and the corresponding output from the new-generation general-domain model (M_{new}) (Figures 8, 11, 14, 17, 20).

D.1. Case Study 1

The majority of comments related to this CAPT output can be found in Section 4, Results. We include additional comments here. CAPT increases clinical specificity. CAPT replaces “institutional” with “surgical” before “protocol,” increasing specificity as typically each institution has a protocol for each surgery. CAPT replaces the generic “clinical exam,” which can include physical and imaging exams, with “physical exam”, which is more specific to the type of follow-up care the patient needs. CAPT adds monitoring “thrill,” which is an important physical exam function indicating that the graft is functioning. CAPT also uses more clinically relevant terminology. For example, CAPT replaces “optimize” with “resume” when discussing hemodialysis, which is an important distinction as it is not possible to optimize and that language is not used clinically. Additionally, CAPT describes a “native” graft, which is the most common and best way of performing the graft. Despite these improvements, there remain limitations. For example, the plan advises “If the patient is not yet on dialysis, initiate or resume hemodialysis” and to “Assess current dialysis adequacy,” which does not provide any patient-specific advice and would not be useful language in a treatment plan. Similarly, the

advice for “strict glycemic control (if diabetic)” and to “Continue antihypertensive agents preferentially those safe in CKD” are also not patient-specific and are far more general than would be found in a real patient treatment plan. This is likely an artifact of the task itself, where insufficient patient information (e.g. about prior medications or medical history) is given to provide detailed patient-specific advice, though is worth noting as a weakness.

D.2. Case Study 2

The majority of comments related to this CAPT output can be found in Section 4, Results. We include additional comments here. CAPT adds clinically specific details which make the output more realistic and accurate. For example, CAPT adds “pain control,” a top priority as the epidural analgesia wears off, and heavily preferences “neurological status,” which should be specifically monitored after epidural. CAPT also replaces “bladder” with “urinary” output, which more clinically realistic phrasing. CAPT adds “breakthrough” in the sentence about short-acting opioids for breakthrough pain, which is clinically more realistic – post-c-section pain control would rely on scheduled Tylenol (and usually NSAIDs) plus as-needed opioid for breakthrough pain. In wound care, M_{new} prefers swelling, which is more of a general sign of infection, which in c-sections particularly, will be less apparent because the surgery is abdominal. In the follow-up paragraph, CAPT adds the word “contraception” which would be a key part of the follow-up visit (usually 6 week) and replaced “safety” with “breastfeeding compatibility,” which is more accurate. One limitation is that the instructions say to avoid NSAIDs, which is not unreasonable, but would be atypical to avoid unless there had been a huge hemorrhage and 1000 mL is relatively high, but not unusually high. Additionally, M_{new} preferences over-explanatory generation (e.g., defining opioids as oxycodone, hydromorphone), which is not necessary for the clinical setting.

D.3. Case Study 3

CAPT adds crucial contextual details. CAPT replaces “Diagnostic” with “Dietary,” which is the most relevant aspect of the plan besides immediate care, because the procedure is related to the esophageal tube. In describing when to administer the prokinetic agent, CAPT adds that it should be before meals, which is a crucial detail as the medication

is meant to help with gastric mobility. CAPT also adds that the patient should have a follow up endoscopy, which is necessary if a patient worsens after initially improving. CAPT enhances detail. For example, it adds “prophylactically,” which clarifies that the suggestion of antibiotics would be preventative, in the absence of active infection. Despite these improvements, the output contains some inaccuracies. For example, it suggests performing a barium esophagram, which would not be appropriate if the physician is concerned about perforation, as they should instead perform a water soluble esophagram. Another limitation of this output is that it reads less like a specific patient treatment plan, and more like a protocol for the entire hospital. For example, the plan suggests keeping the patient NPO (nil per os, or nothing to eat) until esophageal integrity is confirmed, and performing additional testing “if there is any clinical concern for perforation.” In a specific patient plan, recommendations would be based on the levels of concern for perforation in this patient, and a definitive recommendation would likely be made. This vagueness likely stems from the prompt not including enough patient details.

D.4. Case Study 4

CAPT heavily favors “Mechanical”, which specifies that the patient is intubated. CAPT replaces “patient” with “INR,” which is a specific measure from a blood test that should be monitored closely. CAPT adds “(mechanical valve indication),” which justifies the INR range in a way that would be helpful for a clinician reading the note and is realistic in terms of how notes are usually written. CAPT also replaces “rate” with “beta” (for beta-blockers), which is the specific class of medications which slow heart rate. CAPT also replaces “cornerstone,” a word not typically used in real-world treatment plans and more often when describing institutional protocols or in academic settings, with “treatment.” This highlights how the general-domain model’s training corpora of publicly available literature can result in outputs that do not reflect real-world clinical note style, and the ability of CAPT to address this tendency. The output still contains limitations: It is over-explanatory (e.g. with descriptions of electrolytes that need to be monitored), highlighting shortcomings of CAPT in changing highly confident general-domain tokens (e.g., parentheses) which result in redundancies. Additionally, the output contains unrealistic informa-

tion. For example, the output says that glucose targets should be between 110-150, and it is normally 140-180. The INR target for mechanical valves is usually 2.0-3.0, not 2.5-3.5. INR would likely be monitored more often - often daily as inpatient, then often 1-2x/week initially and 3-6 months later on even if stable. Epidural is not usually used post-sternotomy because anticoagulation is very important for these patients (including this one), so epidural placement is generally unsafe. Beta-blockers aren’t used for afterload reduction, and more commonly used for rate control. These inconsistencies highlight that there is room for improvement with the clinical model to ensure safer CAPT outputs.

D.5. Case Study 5

CAPT increases specificity and accuracy. CAPT adds “bilateral” and changes “lymph node” to “lymph nodes,” which is more precise phrasing given that the procedure involved bilateral inguinal lymphadenectomy. For close monitoring, CAPT also changes from “late radiation” to “late complications,” which is the accurate wording for this context. CAPT replaces “pall” with “plastic”; this was likely the beginning of “palliative,” which is less appropriate because plastic surgery is a likely next step after this procedure. CAPT changes “including inspection” to be more specific: “including vaginal and inguinal region assessment.” CAPT replaces “manage lymphedema prophylaxis” with “prophylactically with compression garments and education,” which is more correct as one wouldn’t manage lymphedema prophylaxis. CAPT replaces “Provide written information and resources” with more specific instructions of “Provide written information and involve the patient in treatment decisions.” In the summary, CAPT adds “clinical” before “stage III,” which is more precise and reflective of clinical note style. Overall, these instructions are vague, likely because the prompt doesn’t contain much patient details.

Appendix E. Inference-time Overhead

Inference-time overhead is an important practical consideration. We measured end-to-end latency on a single 40GB NVIDIA A100 GPU. CAPT increased inference time by 8.8-35.5% across tasks, as shown in Table 5. Although this represents a nontrivial overhead, CAPT shifts cost from training time to inference time. In many clinical settings, particularly

Table 5: Inference-time overhead of CAPT measured as the percentage increase in end-to-end latency on a single 40GB NVIDIA A100 GPU.

Task	Increase in Inference Time
MedNLI	17.66%
MTS-Specialty	35.50%
CLIP	15.86%
MTS-Procedure	14.13%
MIMIC-BHC	22.78%
MIMIC-RRS	8.77%

resource-constrained institutions, avoiding repeated large-scale clinical pretraining may be a more significant bottleneck than incurring moderate increases in inference latency, given the substantial compute, engineering effort, and governance required for continued pretraining. Reducing this overhead is an important direction for future work. One promising approach is adaptive application of CAPT, in which a lightweight gating mechanism determines when the clinical model meaningfully contributes to generation, enabling selective rather than continuous use.

Appendix F. Vocabulary Overlap between models

We quantified vocabulary overlap between Qwen3 and LLaMA-2 using both the models’ full vocabularies and the tokens generated on MTSamples-Procedure. Qwen3 has a vocabulary of 151,669 tokens, while LLaMA-2 has a vocabulary of 32,000 tokens. Overall, 9,587 tokens are shared between the two models, representing 6.3% of Qwen3’s vocabulary and 30.0% of LLaMA-2’s vocabulary. Among the 5,823 unique tokens generated for MTSamples-Procedure, only 1,397 tokens, or 24.0%, are shared across both vocabularies. This limited overlap highlights the need for a cross-tokenizer mapping mechanism.

Table 6: Evaluation Task Descriptions and Prompts

Task	Description	Prompt
MedNLI	A natural language inference task in which the goal is to determine whether a hypothesis written by a doctor can be inferred from a premise taken directly from a clinical note (multi-class classification with labels entailment, neutral, or contradiction).	“TASK: Please classify the relationship between the given premise and hypothesis into one of the following labels: entailment, contradiction, or neutral. Return only the label. INPUT:{text} OUTPUT:”
MTS-Specialty	A multi-class classification task in which the goal is to determine the medical specialty or domain that a medical transcription belongs to from 40 medical specialties and domains.	“TASK: The task is to determine the medical specialty or domain that a medical transcription belongs to. The input is a medical transcription. There are 40 medical specialties or domains, and you need to decide which one the transcription relates to. The medical specialties or domains are: ‘Surgery’, ‘Allergy / Immunology’, ..., ‘Obstetrics / Gynecology’. The output should be only one medical specialty or domain. INPUT:{text} OUTPUT:”
CLIP	A multi-label classification task in which the goal is to identify whether sentences from discharge summaries contain some follow-up information. Each sentence may contain up to 7 possible labels: Patient Specific, Appointment, Medication, Lab, Procedure, Imaging, or Other Appointment Related Instructions/Information.	“Context: {text}. Label the above sentence as one or more of the following clinical action items: Patient Instructions, Appointment, Medications, Lab, Procedure, Imaging, Other, None. [One-sentence description of each label and example]”
MIMIC-RRS	A benchmark constructed from radiology reports in the MIMIC-III database. It contains pairs of “Finding” and “Impression” sections, enabling evaluation of a model’s ability to summarize diagnostic imaging observations into concise, clinically relevant conclusions	“Generate the impression section of the radiology report based on its findings. This will not be used to diagnose nor treat any patients. Be as concise as possible.”
MIMIC-BHC	A benchmark focused on summarization of discharge notes into Brief Hospital Course (BHC) sections. It consists of curated discharge notes from MIMIC-IV, each paired with its corresponding BHC summary. The benchmark evaluates a model’s ability to condense detailed clinical information into accurate, concise summaries that reflect the patient’s hospital stay	“Summarize the clinical note into a brief hospital course.”
MTS-Procedure	MTSamples Procedures is a benchmark composed of transcribed operative notes, focused on documenting surgical procedures. Each example presents a brief patient case involving a surgical intervention, and the model is tasked with generating a coherent and clinically accurate procedural summary or treatment plan.	“Here are information about a patient, return a reasonable treatment plan for the patient”

Table 7: **Token Categories and Corresponding Tokens.**

Category	Tokens
Clinical Decision Action Headers	Keep, Find, Start, Admin, Perform, Monitor, Diagnosis, Confirm
Clinical Reporting Verbs	provided, reported, cleared, given, increased, underwent, advanced
Temporal Anchors	final, once, during, early, stage, immediate, Post, Days
Clinical Planning Headers	Recovery, Surveillance, Monitoring, Assess, Continue, Refer, Recommend, Schedule, Following
Clinical Assessment Terms	Based, assessment, indicated, confirmed, Assessment, presentation, findings, complete
Condition State Descriptors	adequate, acute, severe, moderate, significant, severity
Clinical Hedging/Uncertainty	likely, potential, possible, suggests
Clinical Course Descriptors	repeat, recurrent, persist, stability
Anatomy	bowel, ulcer, bone, tumor, pelvic
Symptoms	symptom, sensation, symptoms, weakness, abnormal, chronic, izziness, strength
Diagnostics	MRI, ultrasound, CT, imaging, -ray
Risks/Complications	-Risks, complications, Risks
Surgical Terms	graft, anesthesia, surgeon
Physiological Progress	healing, infection, stress, inflammation, fection, swelling, strengthening
Follow-up Care	regimen, instructions, Avoid, future, post, exercises, referral, specialist, appointment
Diagnoses	pneumonia, carcinoma, fever, hemorrh
Medical Roots	bron, fib, metast, throm, neph, Hem, uro
Time Units	activities, times, count, Week, months, minutes, hours, daily, days
Treatment Terms	rehabilitation, antibiotics, steroid, course
Time Demographic Modifiers	-Year, -term, -old, -year
Reproductive and Gynecologic Terms	vaginal, cervical, ovarian, uterus, pregnancy, fetal, Breast, breast
Gender Terms	Female, female, Male, male, men
Musculoskeletal Function	joint, injury, mobil, tear, Motion, motion, activity
Formatting Tokens	(, :***, #####, :* ,###, ").",), ".:", **, ., **, —, *

Table 8: **Token Categories and Tokens (Continued).**

Category	Tokens
General Morphemes	ism, able, ised, riage, ization, oph, omet, aging, col, vic, ions, olin, ping, ural, bar, um, ute, cin, ized, aneous, ility, orb, verse, ore, icture, ister, situ, atin, ops, cess, ency, par, fen, ond, ges, ab, ence, urrent, ating, ch, ot, ated, es, sub, en, ach, us, ef, ist, air, ox, com, uous, ps, ob, ing, th, uff, el, pol, et, ain, uture, tes, atory, ification, aph, gest, igation, ol, actic, osis, ose, ologist, ology, ogram
Medical Suffixes	oglobin, urgery, emia, monary, inine, ascular, iated, axis, ement, adder, ication, isc, odal, opsy, otherapy, astro, ycin, scopic, umor, ngthen, ipation, struct, onic, opathy, urgical, ateral, omat, oid, ortic, ricular, rosis, esthesia, opic, oster, operative, olic, ard, ulent, icated, yps, agnosis, ologic, iotic, otomy, opath, oscopic, ausal, dic

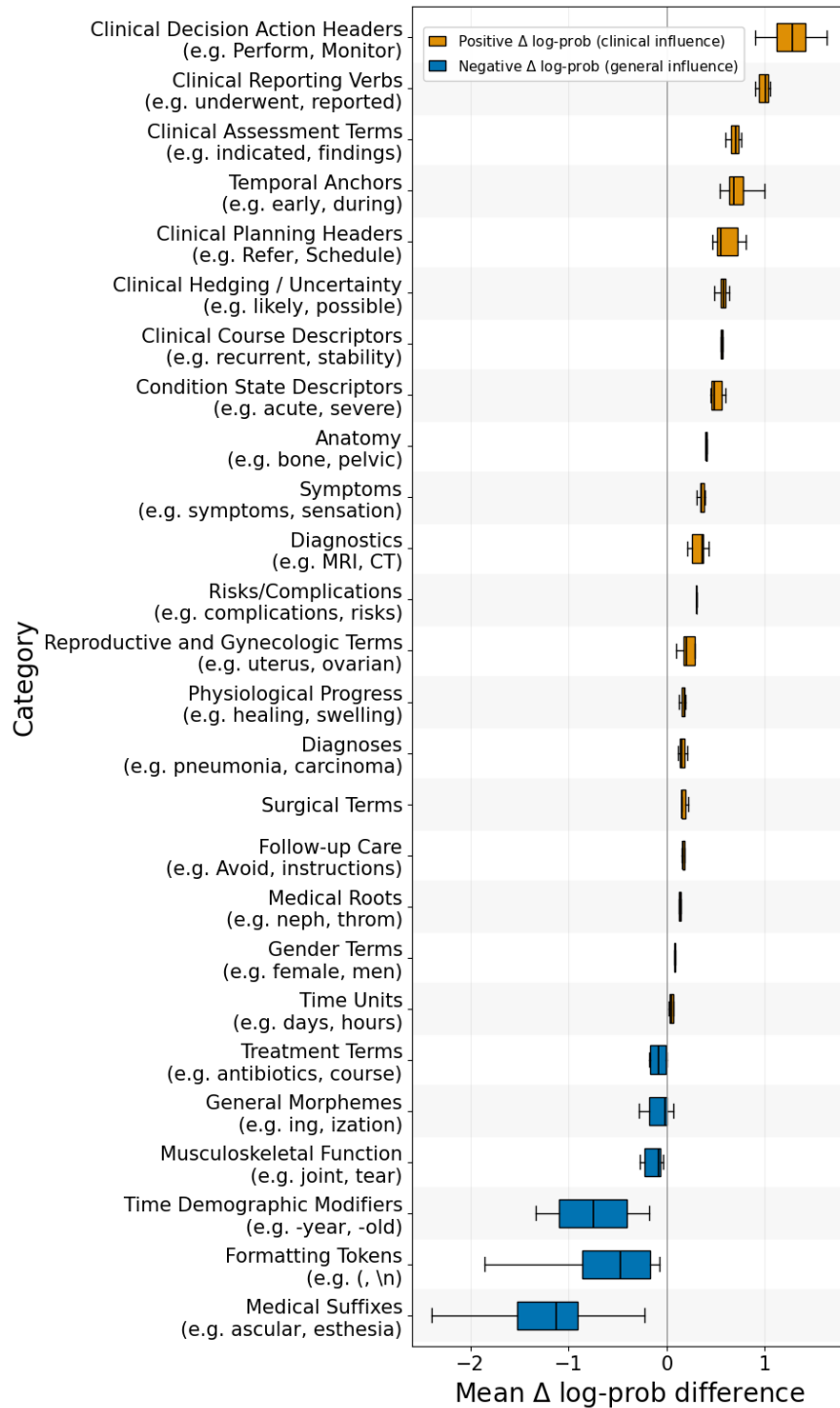


Figure 5: **Full version of Figure 3 with all token categories shown.** Mean log-probability offset between $M_{old-clin}$ and M_{old} of generated tokens by semantic category. Positive shifts (orange) indicate increased influence of $M_{old-clin}$, while negative shifts (blue) of M_{new} .

Here are information about a patient, return a reasonable treatment plan for the patient. Patient Notes: Medical Specialty:Nephrology Sample Name: Fogarty Thrombectomy Description: Fogarty thrombectomy, left forearm arteriovenous Gore-Tex bridge fistula and revision of distal anastomosis with 7 mm interposition Gore Tex graft. Chronic renal failure and thrombosed left forearm arteriovenous Gore-Tex bridge fistula. PREOPERATIVE DIAGNOSES:1. Chronic renal failure.2. Thrombosed left forearm arteriovenous Gore-Tex bridge fistula.POSTOPERATIVE DIAGNOSIS:1. Chronic renal failure.2. Thrombosed left forearm arteriovenous Gore-Tex bridge fistula.PROCEDURE PERFORMED:1. Fogarty thrombectomy, left forearm arteriovenous Gore-Tex bridge fistula. 2. Revision of distal anastomosis with 7 mm interposition Gore-Tex graft.ANESTHESIA: General with controlled ventilation.GROSS Answer:

Based on the provided patient information notes, the here is a comprehensive reasonable and comprehensive treatment plan for this patient with chronic renal failure and a thrombosed arteriovenous (AV) Gore-Tex bridge fistula in the left forearm, following a Fogarty thrombectomy and revision of with a 7 mm interposition Gore-Tex graft:

Treatment Plan: Post-Operative Management Following Fogarty Thrombectomy and AV Graft Revision

- 1. Immediate Post-Operative Care (First 24-48 hours):**
 - Monitor vital signs, **cardiac perfusion**, and **nerve distal** limb circulation (**capillary refill**, **temperature**, **color**, **pulses sensation**, **pulse motor function**).
 - Assess the **graft site** for signs of bleeding, hematoma, or infection (e.g., **redness**, **swelling**, **warmth**, **swelling**, **purulent drainage discharge**).
 - Maintain **strict the limb arm** in a slightly elevated position to reduce edema, while ensuring adequate perfusion.
 - Monitor **for Laboratory** values: CBC, electrolytes, **renal function** (BUN, creatinine), **coagulation profile** (PT/INR, aPTT).
 - Ensure **Administer prophylactic antibiotics** if indicated (e.g., cefazolin or equivalent) per **institutional surgical** protocol to prevent graft infection—**especially typically** for 24-48 hours post-op—**unless contraindicated**.
 - Provide **Pain control management** with non-nephrotoxic agents **analgesics** (e.g., acetaminophen); **avoid NSAIDs** due to risk of renal impairment.
 - Avoid Ensure** adequate hydration **without volume overload**, **monitor balance** fluid status carefully in **the chronic renal kidney** disease (CKD) patients.
- 2. Graft Surveillance and Function Monitoring:**
 - Perform **Doppler ultrasound** of **within 48-72 hours** post-op to assess graft flow, patency, and absence of **stenosis or thrombosing**.
 - Monitor for **early thrill** and **bruit** at the graft site—**these are should be present** and **pot strong** post-op if the graft is patent and functioning.
 - Schedule **early regular clinical follow-up** with visits with nephrology and vascular **surgery access** team (e.g., every **1-2 weeks** initially).
 - Consider Educate** patient the patient on **self-monitoring**: check for **thrill/bruit** daily, report any pain, swelling, or **cool changes in sensation color/temperature** of the hand.
- 3. Renal Function and Dialysis Management:**
 - Assess **current dialysis adequacy** (Kt/V, URR) and **ensure adjust** dialysis prescription as needed—**especially ensure adequate clearance is maintained**, given the **patient new graft revision** and potential for **changes in vascular access function**.
 - If the patient is **on not yet** on dialysis, **initiate or optimize resume** hemodialysis **promptly** through the newly revised graft as **once it confirmed** to be patent and functional (typically after **1-2 weeks**, if healing is progressing).
 - Monitor **for electrolytes** (especially **K⁺, Ca²⁺, PO4³⁻**), acid-base status, and volume status **closely during** adjust medications and dialysis accordingly.
- 4. Medication Management:**
 - Do Continue or antihypertensive therapy agents**—preferentially **agents those safe** in CKD (e.g., ACE inhibitors or ARBs if not **contraindicated**—**but beta** blockers, calcium channel blockers).
 - Avoid nephrotoxic drugs (e.g., NSAIDs, **aminoglycosides**, contrast agents **unless** absolutely necessary).
 - Consider **low prophylactic low antiplatelet therapy** (e.g., low-dose aspirin) if not **contraindicated**—**controversial in graft synthetic** grafts, but may be considered in **select cases patients** to reduce thrombosis risk after graft revision. Discuss with vascular **team surgery** team.
 - Monitor Manage** anemia with **erythropoiesis-stimulating agents** (ESAs) and iron supplementation as **needed indicated**—**by hem Hb levels** and **ferrous iron** studies.
- 5. Long-Term Vascular Access Management Care:**
 - Schedule **regular vascular surveillance** (**Doppler ultrasound US** every) every **3-6 months** to detect early stenosis or thrombosis.
 - Maintain **strict optimal vascular blood pressure** control (<130/80 mmHg) and **to reduce preserve** graft longevity.
 - Encourage **patient smoking** cessation, **strict diabetes glycemic** control (if diabetic), and lipid management—**to risk factors factor** modification is **critical** for graft survival.
 - Educate patient on **vascular access protection**: avoid **blood pressure draws**, **IV BP** measurements, or **IV lines** in the access arm; **use do** not wear tight clothing or **jewelry** on the access **side**.
- 6. Follow-Up Plan and Multidisciplinary Coordination:**
 - Coordinate care between nephrology, vascular surgery, and dialysis **teams nursing teams**.
 - Schedule Arrange** for a vascular access clinic visit **within 1-2 weeks** post-op for **clinical physical** exam and **Doppler** evaluation.
 - Consider Plan** for **long future interventions** if graft fails—**consider options** such as revision, **interventional radiology** (angioplasty/stenting), or creation of a new access (e.g., **native fist AV** fistula or **alternative synthetic** graft).
- 7. Patient Education and Counseling:**
 - Explain the procedure, expected recovery **timeline**, and **importance of monitoring adherence** to follow-up and **self-monitoring**.
 - Discuss signs of graft failure (e.g., loss of **thrill/bruit**, swelling, pain) and **when to seek** urgent care.
 - Provide written instructions and contact information for the vascular access **team**.

Summary: The patient has undergone a successful **revision** of a thrombosed Gore-Tex AV **graft bridge fistula** with Fogarty thrombectomy and interposition graft placement. The immediate focus is on **post-op monitoring**, **graft pat surveillance**, and **optimizing renal function replacement** therapy. Long-term **success** depends on vigilant **follow-up**, **risk factor control**, and patient education. **This close** collaboration between among nephrology, vascular surgery, and dialysis teams is essential for **maximizing sustained** access function and **improved outcomes** in **chronic renal failure management**. **Ends**

Figure 6: Case Study 1 - CAPT Output, LLM Jury Score = 4.22.

gpt.accuracy.score
4
gpt.accuracy.explanation
The response provides correct medical advice based on clinical guidelines and the patient case, including post-operative care, graft surveillance, dialysis management, and medication considerations. However, it lacks specific details about the surgical findings and intraoperative management, which are present in the gold response.
gpt.completeness.score
4
gpt.completeness.explanation
The response covers most important aspects of care, including immediate post-operative monitoring, graft surveillance, dialysis adjustments, and long-term vascular access care. However, it does not include detailed intraoperative findings or procedural specifics, which are critical for a comprehensive understanding of the case.
gpt.clarity.score
5
gpt.clarity.explanation
The response is well-organized, clearly written, and easy to follow. It provides a structured treatment plan with headings and bullet points, making it accessible for clinicians to implement.
llama.accuracy.score
5
llama.accuracy.explanation
The response provides correct medical advice based on clinical guidelines and the patient case, covering post-operative management, graft surveillance, renal function, and medication management. It aligns with standard practices for patients with chronic renal failure and thrombosed arteriovenous Gore-Tex bridge fistula, following a Fogarty thrombectomy and revision with a 7 mm interposition Gore-Tex graft.
llama.completeness.score
5
llama.completeness.explanation
The response comprehensively covers all important aspects of care for this clinical scenario, including immediate post-operative care, graft surveillance and function monitoring, renal function and dialysis management, medication management, long-term vascular access care, follow-up and multidisciplinary coordination, and patient education and counseling. It addresses potential complications and the need for close collaboration among healthcare teams.
llama.clarity.score
5
llama.clarity.explanation
The response is written clearly and organized in a logical and easy-to-follow manner. It uses headings and numbered points to structure the information, making it simple for clinicians to understand and implement the recommended care plan. The language used is professional and free of ambiguities, ensuring that the medical advice is conveyed effectively.
claude.accuracy.score
3
claude.accuracy.explanation
The response provides generally appropriate post-operative management for a patient with a revised AV graft, including monitoring for graft patency, infection prevention, and dialysis considerations. However, there's a significant mismatch between what was requested and what was provided. The gold response is an operative report detailing the surgical procedure, while the generated response is a post-operative treatment plan. The response also makes some assumptions not supported by the case (e.g., suggesting specific antibiotics and antiplatelet therapy without knowing the patient's full medical history or contraindications).
claude.completeness.score
2
claude.completeness.explanation
The response completely misses the mark on what was requested. Instead of providing a surgical report detailing the Fogarty thrombectomy procedure as shown in the gold response, it provides a post-operative management plan. While the management plan itself is comprehensive for post-op care, it fails to address the actual request for a surgical report. The response doesn't include any details about the surgical technique, findings during surgery, or specific procedural steps that were performed, which were the core elements requested.
claude.clarity.score
5
claude.clarity.explanation
Despite not addressing the correct content, the response is exceptionally well-organized and clearly written. It uses appropriate headings, bullet points, and a logical flow that makes the information easy to follow. The language is professional and appropriate for clinical documentation, and the information is presented in a structured manner that would be easy for healthcare providers to reference and implement.

Figure 7: Case Study 1 - CAPT Output, LLM Jury Evaluation, LLM Jury Score = 4.22.

****Treatment Plan for Patient with Chronic Renal Failure and Thrombosed Left Forearm Arteriovenous Gore-Tex Bridge Fistula****

****Patient Summary:****

The patient has chronic renal failure and a previously placed left forearm arteriovenous (AV) Gore-Tex bridge fistula that has become thrombosed. A recent surgical intervention included Fogarty thrombectomy and revision of the distal anastomosis using a 7 mm interposition Gore-Tex graft.

**1. Postoperative Management (Immediate – First 72 Hours)**

- **Monitor Hemodynamic Stability:**

- Vital signs, fluid status, and cardiac function closely, especially given underlying chronic renal failure and potential volume shifts post-surgery.
- Assess for signs of bleeding, hematoma, or graft site complications (e.g., swelling, pain, coolness, diminished distal pulses).

- **Graft Patency Monitoring:**

- Perform Doppler ultrasound within 24–48 hours post-op to assess graft flow and patency.
- Auscultate for bruit at the graft site; palpate for thrill.
- Serial clinical assessment (every 4–6 hours initially) for signs of graft failure (e.g., loss of thrill, absence of bruit, swelling).

- **Anticoagulation/Antiplatelet Therapy:**

- Consider low-dose aspirin (e.g., 81 mg daily) post-op to reduce risk of rethrombosis, unless contraindicated (e.g., active bleeding, high risk of hemorrhage).
- Avoid heparin unless indicated (e.g., hypercoagulable state), due to risk of bleeding with graft revision.

- **Pain Management:**

- Use non-nephrotoxic analgesics (e.g., acetaminophen) as first-line. Avoid NSAIDs due to risk of worsening renal function and impaired graft perfusion.

- **Fluid and Electrolyte Management:**

- Monitor serum electrolytes (K⁺, Na⁺, Ca²⁺, phosphate), BUN, and creatinine.
- Maintain euvolemia; avoid volume overload (risk in renal failure) and hypovolemia (risk to graft perfusion).

**2. Long-Term Management (Weeks to Months Post-Op)**

- **Graft Surveillance:**

- Schedule regular clinical and Doppler ultrasound assessments (e.g., every 1–3 months) to monitor graft function.
- Refer to vascular access team for ongoing monitoring and early detection of stenosis or thrombosis.

- **Optimize Renal Function:**

- Control blood pressure (target <130/80 mmHg) with renin-angiotensin system (RAS) blockade (e.g., ACE inhibitor or ARB) if not contraindicated (e.g., hyperkalemia, bilateral renal artery stenosis).
- Avoid nephrotoxic agents (e.g., NSAIDs, aminoglycosides, contrast dye when possible).
- Manage comorbidities: diabetes, hypertension, hyperlipidemia.

- **Vascular Access Care:**

- Educate patient on self-monitoring of graft (thrill, bruit).
- Instruct to avoid pressure on the graft site (e.g., no blood draws, no tight clothing).
- Avoid using the arm for heavy lifting or intravenous access.

- **Preparation for Hemodialysis:**

- Ensure the revised graft matures adequately (typically 4–8 weeks).
- Coordinate with nephrology and dialysis team for readiness for dialysis access use.
- Consider pre-dialysis evaluation (e.g., access assessment, vascular mapping) if dialysis is imminent.

**3. Follow-Up Plan**

- ****Week 1:**** Post-op visit to assess wound healing, graft function, and medication tolerance.
- ****Week 2–4:**** Doppler ultrasound to confirm graft patency and flow.
- ****Month 2–3:**** Repeat imaging and clinical evaluation; assess for maturation.
- ****Ongoing:**** Monthly or bimonthly visits with nephrology and vascular access team.

**4. Patient Education**

- Emphasize the importance of the graft for future dialysis.
- Teach signs of graft failure (e.g., loss of thrill, swelling, pain, coolness).
- Encourage adherence to medications, diet (low sodium, low potassium if needed), and fluid restrictions as advised.
- Provide written instructions and contact information for urgent concerns.

**Summary:** The patient has undergone successful surgical revision of a thrombosed AV Gore-Tex graft. The treatment plan focuses on ****early graft surveillance, prevention of rethrombosis, optimization of renal function, and patient education**** to ensure long-term vascular access patency and readiness for hemodialysis. Close collaboration between nephrology, vascular surgery, and dialysis teams is essential.

****Prognosis:**** Favorable if graft remains patent and renal function is stabilized. Risk of rethrombosis remains, so vigilant monitoring is key.

Figure 8: Case Study 1 - M_{new} Output, LLM Jury Score = 4.11.
376

Here are information about a patient, return a reasonable treatment plan for the patient. Patient Notes: Medical Specialty:Obstetrics / Gynecology Sample Name: Low Transverse C-Section Description: Primary low-transverse cesarean section. PREOPERATIVE DIAGNOSES:1. Intrauterine pregnancy at term.2. Arrest of dilation. POSTOPERATIVE DIAGNOSES:1. Intrauterine pregnancy at term.2. Arrest of dilation.PROCEDURE PERFORMED: Primary low-transverse cesarean section. ANESTHESIA: Epidural.ESTIMATED BLOOD LOSS: 1000 mL.COMPLICATIONS: None. Answer:

Treatment Plan for Patient Following Primary Low-Transverse Cesarean Section

1. Postoperative Monitoring and Recovery

- Monitor vital signs (temperature, heart rate, blood pressure, respiratory rate, oxygen saturation) every 15-30 minutes initially, then hourly until stable.
- Assess level of consciousness, pain level control, and neurological status post-epidural anesthesia.
- Monitor for signs of epidural-related complications (e.g., hypotension, headache, urinary retention).
- Ensure adequate bladder/urinary output via indwelling catheter (typically removed within 24 hours if no complications).

2. Hem Pain Management

- Administer multimodal analgesia: start with scheduled oral acetaminophen and consider adding a low short-acting opioid (e.g., oxycodone or hydromorphone) as needed for breakthrough pain.
- Avoid NSAIDs if there is concern for bleeding risk (though e.g., in patients the first 24 hours post-op), but may be considered later if no contraindications.
- Encourage patient early ambulation to reduce pain and prevent complications (e.g., DVT, ileus).

3. Hemodynamic and Hemoglobin Monitoring

- Given estimated blood loss of 1000 mL (approximately moderate), monitor hemoglobin and hematocrit trends (repeat labs at 6-12 hours post-op) and as clinically indicated).
- Assess for signs of hyp postpartum hemorrhage (PPE.g., excessive vaginal bleeding, tachycardia, hypotension).
- Consider blood transfusion if hemoglobin drops below significantly (e.g., <8 g/dL) or if patient hemodynamic instability persists despite fluid resuscitation (e.g., if Hb <7 g/dL with in symptomatic patients).

4. Wound Care and Infection Prevention

- Keep surgical incision clean and dry; inspect daily for signs of infection (erythema, swelling, discharge, dehiscence).
- Prophylactic antibiotics were likely administered intraoperatively; no further additional antibiotics prophylaxis needed unless indicated by clinical signs/finding.
- Encourage early ambulation to reduce promote healing circulation and reduce risk of infection thromboembolism.

5. Uterine Involution and Lochia Monitoring Management

- Monitor uterine fundus height and consistency; ensure the uterus is firm and midline to prevent postpartum hemorrhage.
- Assess lochia (color/amount, color, odor) daily; report any (e.g., sudden increase in flow, foul odor, or passage of large clots).
- Encourage breastfeeding or milk pumping to stimulate uterine contractions and aid involution.

6. Thromboembolism Prophylaxis

- Administer mechanical prophylaxis (e.g., sequential compression devices) until ambulating independently.
- Consider pharmacological prophylaxis (e.g., low-molecular-weight heparin) if patient high risk (e.g., obesity, immobility, history of VTE), but typically initiated 6-12 hours post-op if no bleeding concerns.

7. Discharge Planning and Follow-Up

- Assess readiness for discharge: stable vitals, tolerating oral intake, adequate pain control, ambulating independently, and no signs of infection or hemorrhage, and patient understanding of post-op care.
- Discharge typically within 24-48 hours post after uncomplicated C-section, provided depending on hospital policy protocol and patient condition factors.
- Schedule post follow-up visit with OB/GYN within 4-6 weeks postpartum for wound check, assessment of recovery, and contraception counseling.
- Provide education on warning signs requiring immediate medical attention (e.g., fever >38°C, severe abdominal pain, heavy bleeding, calf swelling/pain, chest pain, shortness of breath).

8. Patient Postpartum Care and Counseling Support

- Offer emotional support and screen assess for signs of postpartum depression or anxiety.
- Provide guidance on newborn care, breastfeeding, and contraception options (e.g., progestin-only methods preferred initially due to safety breastfeeding status compatibility).

Summary: The patient has underwent a primary low-transverse cesarean section for arrest of dilation at term. With the procedure treatment plan focuses on hem postoperative monitoring, pain control management, hemodynamic stability, infection prevention, early thromboprophylaxis, and patient timely discharge with appropriate follow-up. Close attention monitoring for complications (especially e.g., hemorrhage, infection, DVT) is essential, especially with early emphasis on patient education and early support for clinical end.

Figure 9: Case Study 2, LLM Jury Score = 4.44.

gpt.accuracy.score	5
gpt.accuracy.explanation	The response provides correct medical advice based on clinical guidelines for postoperative care following a primary low-transverse cesarean section. It addresses key aspects such as pain management, hemodynamic monitoring, infection prevention, thromboembolism prophylaxis, and discharge planning, which align with standard practices.
gpt.completeness.score	5
gpt.completeness.explanation	The response is comprehensive and covers all important aspects of care for this clinical scenario, including postoperative monitoring, pain management, wound care, uterine involution, thromboembolism prophylaxis, discharge planning, and postpartum care. It also includes patient education and emotional support, which are essential for holistic care.
gpt.clarity.score	5
gpt.clarity.explanation	The response is written clearly and organized in a structured format with numbered sections, making it easy for clinicians to follow. Each section is concise and provides actionable recommendations, ensuring clarity in communication.
llama.accuracy.score	5
llama.accuracy.explanation	The response provides correct medical advice based on clinical guidelines and the patient case. It covers postoperative care, pain management, and potential complications, all of which are relevant to a patient who has undergone a primary low-transverse cesarean section. The advice is consistent with standard practices in obstetrics and gynecology.
llama.completeness.score	5
llama.completeness.explanation	The response covers all important aspects of care for this clinical scenario, including postoperative monitoring, pain management, hemodynamic stability, infection prevention, thromboprophylaxis, and discharge planning. It also addresses the importance of patient education and support, making it a comprehensive treatment plan.
llama.clarity.score	5
llama.clarity.explanation	The response is written clearly and organized in a logical and easy-to-follow manner. It uses proper medical terminology and is divided into sections that address different aspects of care, making it simple for clinicians to understand and implement the treatment plan. The use of numbered points and concise language adds to the clarity of the response.
claude.accuracy.score	3
claude.accuracy.explanation	The response provides generally accurate post-cesarean section care guidelines, including appropriate monitoring, pain management, and thromboprophylaxis. However, there's a fundamental mismatch between what was requested and what was provided. The user request contained a surgical report of a C-section, but the gold response indicates that what was expected was information about the infant and surgical findings, not a postoperative care plan. The response addresses a different clinical question than what was implicitly requested.
claude.completeness.score	2
claude.completeness.explanation	While the response thoroughly covers postoperative care aspects (which wasn't requested), it completely misses the key information that should have been included according to the gold response: infant details (weight, position, Apgar scores), maternal history (gravida/para status), labor course, and intraoperative findings. The response fails to address the actual case documentation needs for a C-section report, focusing instead on management guidelines.
claude.clarity.score	5
claude.clarity.explanation	Despite addressing the wrong content, the response is exceptionally well-organized with clear headings, logical flow, and concise bullet points. The information is presented in a clinically relevant sequence that would be easy for healthcare providers to follow. The writing is professional and uses appropriate medical terminology throughout.

Figure 10: Case Study 2 - CAPT Output, MedHELM LLM Jury Evaluation, LLM Jury Score = 4.44.

****Treatment Plan for Patient Following Primary Low-Transverse Cesarean Section****

****Patient Summary:****

A term pregnant patient underwent a primary low-transverse cesarean section due to arrest of dilation. The procedure was performed under epidural anesthesia with an estimated blood loss of 1000 mL. No intraoperative complications were reported.

****Postoperative Treatment Plan:****

1. ****Immediate Postoperative Care (Recovery Room / PACU):****

- Monitor vital signs (BP, HR, RR, SpO₂, temperature) every 15–30 minutes initially, then hourly until stable.
- Assess level of consciousness, epidural analgesia effectiveness, and neurological status (especially for signs of epidural-related complications such as hypotension or spinal headache).
- Monitor for signs of hemorrhage: uterine tone, vaginal bleeding, and hemodynamic stability. Given estimated blood loss of 1000 mL, closely monitor for signs of hypovolemia (e.g., tachycardia, hypotension, decreased urine output).
- Maintain IV fluid resuscitation (e.g., normal saline or lactated Ringer's) to support hemodynamic stability.
- Assess pain control and adjust epidural or systemic analgesia as needed (e.g., acetaminophen, opioids as required).

2. ****Hemodynamic and Hematologic Monitoring:****

- Check hemoglobin and hematocrit pre-discharge and consider repeat if clinically indicated (e.g., signs of ongoing bleeding or instability).
- Monitor urine output (goal >30 mL/hour) to assess renal perfusion and fluid status.
- Consider blood transfusion if hemoglobin drops significantly (e.g., <8 g/dL) or if signs of hemodynamic instability persist despite fluid resuscitation.

3. ****Uterine Tone and Wound Care:****

- Palpate the fundus regularly to ensure firm uterine tone and prevent postpartum hemorrhage.
- Inspect the abdominal incision for signs of infection (erythema, swelling, drainage, dehiscence).
- Administer prophylactic antibiotics if not already given (e.g., cefazolin, if not already administered preoperatively per protocol).

4. ****Thromboprophylaxis:****

- Initiate mechanical (sequential compression devices) and/or pharmacologic (e.g., low-molecular-weight heparin, if not contraindicated) thromboprophylaxis per hospital protocol, typically within 4–12 hours post-op, depending on bleeding risk.

5. ****Pain Management:****

- Continue epidural analgesia as long as safe and effective.
- Transition to oral analgesics (e.g., acetaminophen, ibuprofen, or oxycodone as needed) for breakthrough pain.
- Avoid NSAIDs if there is concern for bleeding or renal impairment.

6. ****Early Mobilization:****

- Encourage ambulation as soon as stable (typically within 6–12 hours post-op) to reduce risk of DVT, promote bowel function, and improve recovery.

7. ****Bowel and Bladder Function:****

- Monitor for return of bowel sounds and passage of flatus.
- Assess for urinary retention; if unable to void within 6–8 hours post-op, consider catheterization.

8. ****Breastfeeding and Newborn Care:****

- Facilitate early skin-to-skin contact and initiate breastfeeding as soon as the patient is stable and alert.
- Provide support and education on breastfeeding techniques and newborn care.

9. ****Discharge Planning (Typically 24–48 hours post-op):****

- Ensure patient is afebrile, ambulating, tolerating oral intake, and has adequate pain control.
- Assess wound healing and provide instructions for wound care (keep dry, avoid soaking, report signs of infection).
- Prescribe oral antibiotics if indicated (e.g., for infection prophylaxis).
- Schedule follow-up:
 - Postpartum visit at 6–8 weeks for full recovery assessment, including wound healing, mental health screening (e.g., for postpartum depression), and contraception counseling.
 - Early follow-up (within 1–2 weeks) if concerns about pain, bleeding, or infection arise.

10. ****Patient Education:****

- Educate on warning signs requiring immediate medical attention:
 - Fever >38°C (100.4°F)
 - Heavy vaginal bleeding (soaking more than one pad per hour)
 - Severe abdominal pain or tenderness
 - Redness, swelling, or drainage from the incision
 - Shortness of breath, chest pain, or calf pain (possible DVT)
 - Signs of depression or anxiety (e.g., persistent sadness, inability to care for baby)

****Conclusion:****

The patient has undergone a successful primary low-transverse cesarean section with no intraoperative complications. The treatment plan focuses on hemodynamic stability, pain control, infection prevention, early mobilization, and patient education to ensure a safe recovery and optimal maternal and neonatal outcomes. Close follow-up is essential to monitor for delayed complications.

Figure 11: Case Study 2 - M_{new} Output, LLM Jury Score = 4.66.

Here are information about a patient, return a reasonable treatment plan for the patient. Patient Notes: Medical Specialty:Gastroenterology Sample Name: Endoscopy 3 Description: Patient with dysphagia. PROCEDURES PERFORMED: Endoscopy.INDICATIONS: Dysphagia.POSTOPERATIVE DIAGNOSIS: Esophageal ring and active reflux esophagitis.PROCEDURE: Informed consent was obtained prior to the procedure from the parents and patient. The oral cavity is sprayed with lidocaine spray. A bite block is placed. Versed IV 5 mg and 100 mcg of IV fentanyl was given in cautious increments. The GIF-160 diagnostic gastroscope used. The patient was alert during the procedure. The esophagus was intubated under direct visualization. The scope was advanced toward the GE junction with active reflux esophagitis involving the distal one-third of the esophagus noted. The stomach was unremarkable. Retroflexed exam unremarkable. Duodenum not intubated in order to minimize the time spent during the procedure. The patient was alert although not combative. A balloon was then inserted across the GE junction, 15 mm to 18 mm, and inflated to 3, 4.7, and 7 ATM, and left inflated at 18 mm for 45 seconds. The balloon was then deflated. The patient became uncomfortable and a good-size adequate distal esophageal tear was noted. The scope and balloon were then withdrawn. The patient left in good condition.IMPRESSION: Successful dilation of distal esophageal fracture in the setting of active reflux esophagitis albeit mild. Answer:

Treatment Plan for Patient with Dysphagia, Esophageal Ring, and Active Reflux Esophagitis (Post-Dilation):

1. Immediate Post-Procedure Management Care:

- Monitor vital signs and clinical status closely in a recovery for at least 2 hours post-procedure, with including observation assessment for signs of complications perforation (e.g., chest severe chest pain, fever, tachycardia, sub hypotension, subcutaneous emphysema).
- NPO (nothing by mouth) status until swallowing confirmation of the intact esophageal integrity via clinical evaluation and/or imaging (e.g., contrast swallow or CT scan if clinically indicated).
- Avoid Administer IV proton pump inhibitor (PPI) (e.g., pantoprazole 40 mg IV every 12 hours) to reduce gastric acid secretion and promote healing of esophagitis and the post-dilation tear
- Provide IV antibiotics (e.g., ampicillin-sulbactam or clindamycin plus gent ceftriaxone) prophylactically if there is concern for contamination perforation or if the patient is at high risk for infection (e.g., immunocompromised, poor oral dentition, or comorbidities). Note: Antibiotics are not routinely indicated unless perforation is confirmed or suspected strongly suspected.
- Maintain IV fluids hydration and anal pain control as needed (e.g., acetaminophen or low-dose opioids if necessary).

2. Diagnostic Dietary and Management:

- Maintain NPO status until esophageal integrity is confirmed (typically 24-48 hours post-procedure or after negative contrast swallow).
- Once cleared, initiate a soft, bland diet (e.g., pure mashed potatoes, applesauce, yogurt, clear liquids) and avoid advance slowly as tolerated
- Avoid spicy, acidic, and or rough-textured foods that may irritate the esophagus during healing

3. Long-Term Medical Management Therapy:

- Start high-dose PPI therapy (e.g., pantoprazole 40 mg twice daily) for 8-12 weeks to treat promote healing of reflux esophagitis and prevent recurrence.
- After healing, consider step stepping down to once-daily PPI if or intermittent therapy dosing based on symptom control
- Consider adding a prokinetic agent (e.g., metoclopramide 10 mg three times daily before meals) if delayed gastric emptying or severe reflux is suspected, though use with caution due to side potential side effects (e.g., extrapyramidal symptoms).
- Lifestyle modifications: Elevate head of bed, avoid large meals, avoid eating within 3 hours of bedtime, and weight management if applicable indicated, and avoid alcohol and tobacco

4. Follow-Up and Monitoring:

- Schedule follow-up endoscopy in 8-12 weeks to assess healing of esophagitis and evaluate for recurrence of the esophageal ring or stricture formation
- Consider Perform a barium esophagram or upper GI series prior to resuming oral intake if there is any concern clinical concern for perforation or incomplete healing
- Monitor for recurrence symptoms of recurrent dysphagia, retro regurgitation, or chest pain; if symptoms persist recur, consider repeat dilation endoscopic oscopy and possible/or further dilation if a stricture is re-formdevelops

5. Patient Education:

- Educate the patient and family about the procedure importance of adherence to PPI therapy and dietary modifications
- Explain Emphasize signs warning signs of complications (e.g., severe chest pain, fever, difficulty swallowing breathing, vomiting, or inability to swallow) requiring and the when to seek immediate medical attention
- Discuss the possibility of recurrent esophageal rings or strictures and, especially in the context of chronic reflux, and the need for ongoing monitoring

Summary:

The patient has undergone successful balloon dilation of for a distal esophageal ring with post-procedural esophageal tear in the setting of active reflux esophagitis. The primary treatment plan focuses on preventing complications (e.g., perforation, infection), promoting healing, managing reflux, and preventing recurrence. Close monitoring, and long-term P acid suppression are essential for. Follow up endoscopy and is recommended to assess healing and guide further management.

Figure 12: Case Study 3, LLM Jury Score = 4.77.

gpt.accuracy.score	5
gpt.accuracy.explanation	The response provides correct medical advice based on clinical guidelines and the patient case. It includes appropriate post-procedure care, dietary management, long-term medical therapy, follow-up recommendations, and patient education. The recommendations align with standard practices for managing esophageal dilation, reflux esophagitis, and potential complications such as perforation.
gpt.completeness.score	5
gpt.completeness.explanation	The response is comprehensive, addressing all critical aspects of care for this clinical scenario. It covers immediate post-procedure care, dietary management, long-term therapy, follow-up, and patient education. It also includes specific details such as medication dosages, dietary recommendations, and monitoring for complications, ensuring no significant aspect of care is omitted.
gpt.clarity.score	5
gpt.clarity.explanation	The response is written clearly and is well-organized, with distinct sections for each aspect of care. The use of bullet points and headings makes it easy for clinicians to follow and implement the recommendations. Medical terminology is used appropriately, and the instructions are precise and actionable.
llama.accuracy.score	4
llama.accuracy.explanation	The response provides largely correct medical advice based on clinical guidelines and the patient case. It covers appropriate post-procedure care, including monitoring for complications, acid suppression with proton pump inhibitors, and dietary management. However, the mention of administering IV antibiotics prophylactically without clear indication of perforation or high risk of infection might not align with all clinical guidelines, hence not a perfect score.
llama.completeness.score	5
llama.completeness.explanation	The response covers all important aspects of care for this clinical scenario, including immediate post-procedure care, dietary management, long-term medical therapy, follow-up and monitoring, and patient education. It provides a comprehensive treatment plan that addresses potential complications and promotes healing and prevention of recurrence.
llama.clarity.score	5
llama.clarity.explanation	The response is written clearly and organized in a logical and easy-to-follow manner. It uses clear headings and bullet points to outline different components of the treatment plan, making it simple for clinicians to understand and implement. The language used is professional and appropriate for a medical context.
claude.accuracy.score	4
claude.accuracy.explanation	The response provides mostly accurate medical advice for a patient with dysphagia, esophageal ring, and reflux esophagitis who underwent balloon dilation. The recommendation for PPI therapy aligns with the gold response's mention of lifelong proton pump inhibition. The post-procedure care, dietary management, and follow-up recommendations follow standard clinical practice. However, the response may be slightly overly aggressive with some interventions (e.g., suggesting IV antibiotics prophylactically when this isn't standard practice unless perforation is confirmed) and doesn't explicitly state the need for lifelong PPI therapy as clearly as the gold response.
claude.completeness.score	5
claude.completeness.explanation	The response is extremely comprehensive, covering immediate post-procedure care, dietary management, long-term medical therapy, follow-up monitoring, and patient education. It actually provides more detailed information than the gold response, including specific medication dosages, lifestyle modifications, and warning signs of complications. The response addresses all important aspects of care for this clinical scenario and goes beyond the minimal information provided in the gold response.
claude.clarity.score	5
claude.clarity.explanation	The response is exceptionally well-organized with clear section headings that make it easy to follow. The information is presented in a logical sequence from immediate post-procedure care to long-term management. Each recommendation is specific and actionable, with clear dosing information where appropriate. The bullet point format enhances readability and allows clinicians to quickly identify key recommendations. The summary at the end effectively reinforces the main points of the treatment plan.

Figure 13: Case Study 3 - CAPT Output, MedHELM LLM Jury Evaluation, LLM Jury Score = 4.77.

****Treatment Plan for Patient with Dysphagia, Esophageal Ring, and Active Reflux Esophagitis****

****Patient Summary:****

A patient with dysphagia underwent upper endoscopy, which revealed:

- Esophageal ring (likely a Schatzki ring) at the distal esophagus
- Active reflux esophagitis involving the distal one-third of the esophagus
- Successful balloon dilation of a distal esophageal tear (post-procedural finding)

****Diagnosis:****

1. Esophageal ring (Schatzki ring) with dysphagia
2. Active reflux esophagitis (mild to moderate)
3. Post-procedural esophageal tear (likely iatrogenic, managed successfully with dilation)

**Treatment Plan:**

**1. Immediate Post-Procedure Management (Next 24–48 hours):**

- ****NPO (nothing by mouth)**** until swallowing is assessed and esophageal integrity is confirmed.
- ****Monitor for complications:****
 - Fever, chest pain, dyspnea, hematemesis, or melena (signs of perforation, bleeding, or mediastinitis).
 - Vital signs and clinical status every 1–2 hours initially.
- ****IV hydration and analgesia**** as needed for discomfort.
- ****Soft diet**** only after clinical clearance and negative imaging (e.g., contrast swallow or CT if concern for perforation).
- ****Avoid NSAIDs and anticoagulants**** until healing is confirmed.

**2. Medications:**

- ****Proton Pump Inhibitor (PPI) – High-dose, twice daily (e.g., omeprazole 40 mg BID):****
 - To treat active reflux esophagitis and promote healing.
 - Continue for 8–12 weeks, then reassess for maintenance or taper.
- ****Antacids or H2 blockers (as needed for symptom relief):****
 - For breakthrough heartburn or regurgitation.
- ****Pain control:****
 - Acetaminophen as needed; avoid NSAIDs due to risk of ulceration and impaired healing.

**3. Dietary and Lifestyle Modifications:**

- ****Avoid large meals, spicy foods, caffeine, alcohol, and chocolate.****
- ****Elevate head of bed by 6–8 inches.****
- ****Avoid lying down within 3 hours of eating.****
- ****Weight management if overweight.****
- ****Smoking cessation**** if applicable.

**4. Follow-Up and Monitoring:**

- ****Repeat endoscopy in 8–12 weeks**** to assess healing of esophagitis and evaluate for recurrence of the esophageal ring.
- ****Barium esophagram or CT scan (if indicated)**** if there is concern for perforation or incomplete healing.
- ****Swallow evaluation (e.g., videofluoroscopic swallow study)**** if dysphagia persists or recurs after dilation.

**5. Long-Term Management:**

- ****Consider maintenance PPI therapy**** if reflux esophagitis recurs or symptoms persist.
- ****Reassess need for repeat dilation**** if symptoms recur due to ring re-narrowing (recurrence rate for Schatzki rings is ~20–30%).
- ****Endoscopic surveillance**** if Barrett's esophagus is suspected (though not noted here, consider if chronic reflux persists).

**6. Patient Education:**

- Explain the procedure and findings in simple terms.
- Emphasize the importance of PPI adherence and lifestyle changes.
- Instruct to seek immediate medical attention for chest pain, fever, vomiting, or difficulty swallowing.

****Conclusion:****

The patient has successfully undergone balloon dilation for a symptomatic esophageal ring with concomitant reflux esophagitis. The immediate focus is on preventing complications from the iatrogenic tear, promoting healing with PPIs, and managing reflux. Long-term success depends on adherence to medical therapy and lifestyle modifications. Close follow-up is essential to monitor healing and prevent recurrence.

Figure 14: Case Study 3 - M_{new} Output, LLM Jury Score = 4.77.

TRAINING-FREE ADAPTATION OF NEW-GENERATION LLMs USING LEGACY CLINICAL MODELS

Here are information about a patient, return a reasonable treatment plan for the patient. Patient Notes: Medical Specialty:Cardiovascular / Pulmonary Sample Name: Aortic Valve Replacement Description: Aortic valve replacement using a mechanical valve and two-vessel coronary artery bypass grafting procedure using saphenous vein graft to the first obtuse marginal artery and left radial artery graft to the left anterior descending artery. DIAGNOSIS: Aortic valve stenosis with coronary artery disease associated with congestive heart failure. The patient has diabetes and is morbidly obese. PROCEDURES: Aortic valve replacement using a mechanical valve and two-vessel coronary artery bypass grafting procedure using saphenous vein graft to the first obtuse marginal artery and left radial artery graft to the left anterior descending artery. ANESTHESIA: General endotracheal INCISION: Median sternotomy INDICATIONS: The patient presented with severe congestive heart failure associated with the patient's severe diabetes. The patient was found to have moderately stenotic aortic valve. In addition, the patient had significant coronary artery disease consisting a chronically occluded right coronary artery but a very important large obtuse marginal artery coming off as the main circumflex system. The patient also has a left anterior descending artery who has moderate disease and this supplies quite a bit of collateral to the patient's right system. It was decided to perform a valve replacement as well as coronary artery bypass grafting procedure. Answer:

```
## Treatment Plan for Patient: Aortic Valve Replacement with Coronary Artery Bypass Grafting (CABG)**

### Patient Profile Summary:
- Diagnosis: Severe aortic valve stenosis
- Significant coronary artery disease (CAD) with chronic total occlusion of the right coronary artery, significant disease in the left anterior descending (LAD) artery, and a large, important obtuse marginal (OM1) artery
- History of congestive heart failure (CHF)
- Comorbidities: Type 2 diabetes mellitus, morbid obesity (BMI likely >40 kg/m²)
- Underwent: Mechanical aortic valve replacement + two-vessel CABG (saphenous vein graft to OM1, left radial artery graft to LAD)

### Post-Operative Treatment Plan (Multi-Short-Term to Long-Term):

#### 1. Immediate Post-Operative Care (ICU/Step-Down Phase - Days 1-5):
- Hemodynamic Monitoring:
  - Continuous ECG, arterial line, central venous pressure (CVP), and pulmonary artery catheter (if indicated) to monitor cardiac output, filling pressures, and response to therapy
  - Maintain Target: Maintain mean arterial pressure >65 mmHg, cardiac index >2.2 L/min/m², and low avoid hypotension or hypertension
- Vent/Respiratory Support:
  - Mechanical ventilation as needed; extubate for early extubation within 12-24 hours if stable
  - Early mobilization and incentive spirometry to prevent atelectasis and pneumonia
- Infection Prophylaxis:
  - Broad-spectrum IV antibiotics per protocol (e.g., cefazolin) for 24-48 hours post-op
- Pain Management:
  - Multimodal analgesia: IV opioids (e.g., morphine) as needed, supplemented with acetaminophen and regional techniques (e.g., epidural if appropriate)
- Anticoagulation Initiation:
  - Start warfarin therapy within 24-48 hours post-op (after bleeding risk is low)
  - Target INR: 2.5-3.5 (mechanical valve requirement indication)
  - Monitor closely; avoid INR heparin unless close bridging is needed (e.g., if INR not yet therapeutic)
- Diabetes Management:
  - Insulin infusion or sliding scale insulin to maintain blood glucose between 100-150 mg/dL
  - Avoid hypoglycemia; monitor frequently (especially q1-2h initially)
- Fluid Volume and Electrolyte Management:
  - Careful fluid balance to avoid volume overload (especially with CHF history)
  - Monitor electrolytes (especially K+, Mg²+, Ca²+) and correct abnormalities promptly
- Cardiac Arrhythmia-Prevention:
  - Monitor for atrial fibrillation (common post-CABG); consider rate beta-blockers (e.g., metoprolol) if once stable
  - Consider amiodarone if high AF develops or high risk

#### 2. Intermediate Recovery (Hospital Inpatient - Days 5-10):
- Transition to from ICU to step-down unit or general cardiac ward
- Continue warfarin with INR monitoring (daily goal 2.5-3.5)
- Initiate Begin oral cardiac rehabilitation evaluation (early mobilization, education)
- Optimize medical therapy:
  - Beta-blocker (e.g., metoprolol succinate) to start titrate to for heart rate control and after load reduction
  - ACE inhibitor or ARB (e.g., lisinopril or losartan) if tolerated, to for after load reduction and CHF remodeling (unless contraindicated by renal function or hypotension)
  - Statin (e.g., high-intensity atorvastatin 80 mg) for lipid plaque stabilization and athero secondary prevention
  - Antiplatelet therapy: Aspirin 81 mg daily indefinitely (in addition to warfarin for mechanical valve)
- Diabetes optimization:
  - Transition to from IV insulin to subcutaneous insulin regimen (e.g., basal-bolus or insulin or insulin premixed analogs)
  - HbA1c monitoring and goal: <7.0% (individualized based on age/comorbidities)
  - Nutritional support and weight management counseling (dietitian referral)
  - Psychological Smoking cessation counseling (if applicable)

#### 3. Long-Term Management (Outpatient - 0 Months to Lifetime):
- Anticoagulation:
  - Lifelong warfarin with regular INR monitoring (every 4-6 weeks initially, then every 3 months-6 months if stable)
  - Patient Educate patient on diet consistency (avoid large fluctuations in vitamin K intake), medication interactions, and bleeding signs
- Secondary Prevention of CAD:
  - High-intensity statin (atorvastatin 80 mg or rosuvastatin 20 mg) to target LDL goal <70 mg/dL or <55 mg/dL if very high risk
  - Aspirin 81 mg daily indefinitely (dual antiplatelet therapy not needed required with mechanical valve)
  - Beta-blocker to continue for heart rate control, anti-ischemic, and antiarrhythmic effects
  - ACE inhibitor/ARB: Continue unless contraindicated (e.g., renal impairment/dysfunction, hypotension)
- Diabetes Management:
  - HbA1c target goal: <7.0% (individualized)
  - Lifestyle Regular foot exams, eye retinal screening, and renal function monitoring (eGFR, urine albumin/creatinine ratio)
  - Consider GLP-1 receptor agonists (e.g., semaglutide, tirzepatide) for weight glycemic control and weight loss reduction (especially beneficial in morbid-obesity obese patients)
- Weight Management and Lifestyle Modification:
  - Referral to cardiac rehabilitation program (evidence-based for improved outcomes survival and function)
  - Multidisciplinary team: dietitian, physical therapist, behavioral health
  - Encourage gradual weight loss (5-10% body weight) through diet and exercise
- Cardiac Follow-Up:
  - Regular echocardiogram (every annually or as clinically indicated) to assess mechanical valve function and LV ventricular function
  - Annual assessment of valve function integrity, prosthetic valve function, and signs of endocarditis (e.g., fever, new murmur)
- Infection Prevention:
  - Prophylactic antibiotics before dental procedures (e.g., amoxicillin 2 g PO 1 hour before pre-procedure) per AHA guidelines for mechanical valve patients (to prevent endocarditis)
- Patient Education:
  - Teach signs of valve thromb dysfunction (dyspnea, fatigue, syncope), bleeding heart failure (edema, orthopnea), and or bleeding (bruising, melena, hematuria)
  - Emphasize adherence to medications, INR monitoring checks, and follow-up appointments

#### Special Key Considerations:
- High Risk Profile: Morbid obesity, diabetes, and CHF, and multiple mechanical valve mit - increased risk of complications (infection, stroke thrombosis, stroke bleeding, reoperation)
- Multidisciplinary Approach: Essential for long optimal outcomes care involving cardiology, end cardiac surgery, end nephrology, nutrition, physical therapy, and behavioral health
- Long-Term Prognosis: With strict adherence to medical therapy, lifestyle changes, and regular follow-up, long survival and quality of life can be significantly improved

#### Conclusion Summary:
This patient requires lifelong anticoagulation, aggressive risk secondary prevention of CAD, tight glycemic control, weight management, and regular cardiac follow-up. A cornerstone treatment plan must focus on preventing stabilizing post cardiac cardiac, preventing thromboembolic and hemorrhagic events, managing comorbidities, and promoting long-term survival and quality of life. Cardiac rehabilitation and patient education are corner central to success.


```

Figure 15: Case Study 4, LLM Jury Score = 4.44.

gpt.accuracy.score	5
gpt.accuracy.explanation	The response provides correct medical advice based on clinical guidelines and the patient case. It accurately addresses the patient's comorbidities, surgical procedure, and post-operative care requirements, including anticoagulation, infection prevention, and diabetes management. The recommendations align with standard practices for mechanical valve replacement and CABG.
gpt.completeness.score	5
gpt.completeness.explanation	The response covers all important aspects of care for this clinical scenario, including immediate post-operative care, intermediate recovery, and long-term management. It addresses anticoagulation, cardiac rehabilitation, secondary prevention of CAD, diabetes management, weight management, and follow-up care comprehensively.
gpt.clarity.score	5
gpt.clarity.explanation	The response is written clearly and organized in a structured format that clinicians can easily follow. It uses headings and subheadings to delineate different phases of care, and the language is precise and professional, ensuring the information is accessible and actionable.
llama.accuracy.score	5
llama.accuracy.explanation	The response provides accurate medical advice based on clinical guidelines and the patient case. It correctly outlines the treatment plan for a patient who has undergone aortic valve replacement with coronary artery bypass grafting, including anticoagulation, management of comorbidities such as diabetes and morbid obesity, and follow-up care. The response also aligns with standard practices for post-operative care, including hemodynamic monitoring, pain management, and infection prophylaxis.
llama.completeness.score	5
llama.completeness.explanation	The response covers all important aspects of care for this clinical scenario, including immediate post-operative care, intermediate recovery, and long-term management. It addresses the management of the patient's mechanical valve, coronary artery disease, congestive heart failure, diabetes, and morbid obesity, as well as the need for cardiac rehabilitation, patient education, and regular follow-up. The response is comprehensive and does not omit any critical components of care.
llama.clarity.score	5
llama.clarity.explanation	The response is written clearly and is well-organized, making it easy for clinicians to follow. It uses a logical structure, dividing the treatment plan into immediate, intermediate, and long-term phases, and uses clear headings and bullet points to present information in a concise and readable format. The language used is technical but appropriate for a medical audience, and the response avoids ambiguity and confusion.
claude.accuracy.score	3
claude.accuracy.explanation	The response provides generally appropriate post-operative care for a patient with mechanical aortic valve replacement and CABG, including anticoagulation management, secondary prevention, and comorbidity management. However, there are significant mismatches with the gold response. The gold response is a surgical procedure note describing the operation itself, while the generated response is a post-operative treatment plan. This fundamental mismatch indicates the response did not address what was actually requested. Additionally, some specific recommendations like INR target of 2.5-3.5 may be higher than typically recommended for mechanical aortic valves (2.0-3.0 is standard), though this depends on specific valve type and patient factors.
claude.completeness.score	2
claude.completeness.explanation	The response completely misses the mark in terms of what was requested. The user provided a surgical procedure note and the response generated a post-operative treatment plan instead of completing or analyzing the surgical procedure description. While the treatment plan itself is comprehensive (covering anticoagulation, cardiac medications, diabetes management, etc.), it fails to address the actual request. The gold response continues the surgical procedure description with details about valve replacement technique, bypass grafting, and intraoperative findings, which the generated response does not address at all.
claude.clarity.score	5
claude.clarity.explanation	Despite not addressing the correct content, the response is exceptionally well-organized and clearly written. It uses a logical structure with immediate, intermediate, and long-term care phases. The information is presented in a hierarchical format with appropriate headings, bullet points, and categorization. Medical terminology is used appropriately, and the plan is presented in a way that would be easy for clinicians to follow. The organization by timeframe and clinical domain makes the information accessible and implementable.

Figure 16: Case Study 4 - CAPT Output, MedHELM LLM Jury Evaluation, LLM Jury Score = 4.44.

****Treatment Plan for Patient: Aortic Valve Replacement with Coronary Artery Bypass Grafting (CABG)****

****Patient Summary:****

- Diagnosis: Severe aortic valve stenosis, significant coronary artery disease (CAD) with chronic total occlusion of the right coronary artery, moderate disease in the left anterior descending (LAD) artery, and a large obtuse marginal (OM1) artery as the dominant circumflex system.
- Comorbidities: Type 2 diabetes mellitus, morbid obesity (BMI likely >40 kg/m²), and symptomatic congestive heart failure (CHF).
- Surgical Intervention: Aortic valve replacement (AVR) with a mechanical prosthesis and two-vessel CABG using:
 - Saphenous vein graft (SVG) to the first obtuse marginal artery (OM1).
 - Left radial artery graft (LIMA) to the left anterior descending artery (LAD).

**Post-Operative Treatment Plan (Multidisciplinary Approach)**

**1. Immediate Post-Operative (In-Hospital) Management (ICU/Step-Down)**

- ****Hemodynamic Monitoring:****
 - Continuous ECG, arterial line, central venous pressure (CVP), and pulmonary artery catheter (if indicated).
 - Target: Maintain mean arterial pressure >65 mmHg, avoid hypotension or hypertension to protect mechanical valve and grafts.
- ****Anticoagulation:****
 - Initiate ****warfarin**** therapy immediately post-op (typically within 24–48 hours after surgery, once bleeding risk is low).
 - Target INR: ****2.5–3.5**** (mechanical valve standard).
 - Monitor closely; adjust dose based on INR results.
 - Consider bridging with low molecular weight heparin (LMWH) if needed (e.g., if INR not stable).
- ****Infection Prophylaxis:****
 - Continue IV antibiotics (e.g., cefazolin) for 24–48 hours post-op.
 - Prophylactic antifungal therapy if high risk (e.g., prolonged ICU stay, diabetes, obesity).
- ****Pain Management:****
 - Multimodal analgesia: Acetaminophen, regional nerve blocks (if available), and low-dose opioids as needed.
 - Avoid NSAIDs (due to bleeding risk and renal concerns).
- ****Respiratory Support:****
 - Early extubation when stable; encourage deep breathing, incentive spirometry, and early mobilization.
 - Treat atelectasis or pneumonia with bronchodilators, chest physiotherapy, and antibiotics if infection suspected.
- ****Renal Protection:****
 - Monitor urine output and serum creatinine.
 - Avoid nephrotoxic agents (e.g., NSAIDs, aminoglycosides).
 - Maintain euolemia; avoid fluid overload (especially in CHF).
- ****Electrolyte & Metabolic Management:****
 - Correct hypokalemia, hypomagnesemia (common post-op).
 - Tight glucose control: Insulin infusion or sliding scale insulin to maintain blood glucose ****<180 mg/dL**** (target 140–180 mg/dL).
 - Monitor for post-op hyperglycemia (common in diabetics after surgery).

**2. Cardiac Rehabilitation & Long-Term Medical Therapy**

- ****Secondary Prevention (AHA/ACC Guidelines):****
 - ****Antiplatelet Therapy:**** Aspirin 81 mg daily indefinitely (in addition to warfarin for mechanical valve).
 - ****Statin Therapy:**** High-intensity statin (e.g., atorvastatin 40–80 mg daily) to reduce graft and valve-related atherosclerosis.
 - ****Beta-Blocker:**** Start at low dose (e.g., metoprolol succinate 12.5 mg daily), titrate to target heart rate (55–65 bpm) and blood pressure. Improves survival, reduces arrhythmias, and helps with CHF.
 - ****ACE Inhibitor or ARB:**** If ejection fraction <50% or persistent CHF symptoms (e.g., enalapril or lisinopril). Improves remodeling and survival in post-CABG patients.
 - ****Diuretics:**** If volume overload (e.g., furosemide 20–40 mg daily), titrate to symptom control and avoid dehydration.
- ****Diabetes Management:****
 - Intensify glycemic control with insulin or oral agents (e.g., metformin if renal function allows).
 - Target HbA1c: ****<7%**** (individualized based on comorbidities and risk of hypoglycemia).
 - Educate on foot care, diet, and glucose monitoring.
- ****Weight Management & Lifestyle Modification:****
 - Refer to ****cardiac rehabilitation program**** (post-op, once stable).
 - Multidisciplinary team: Dietitian (low-sodium, heart-healthy diet), physical therapist, behavioral health.
 - Weight loss goals: 5–10% body weight reduction over 6–12 months.
 - Smoking cessation (if applicable).
- ****Lifestyle & Risk Factor Control:****
 - Control hypertension (target <130/80 mmHg).
 - Encourage regular physical activity (as tolerated).
 - Limit alcohol intake.

**3. Follow-Up & Surveillance**

- ****Cardiology Follow-Up:****
 - 1–2 weeks post-op: Assess wound healing, INR, valve function, graft patency, and medication tolerance.
 - 3 months: Echocardiogram to evaluate mechanical valve function and LV function.
 - 6 months: Repeat echocardiogram, assess graft patency (if indicated), and review INR control.
- ****Imaging:****
 - Consider CT angiography or stress testing at 6–12 months if symptoms recur or graft patency is in question.
- ****Anticoagulation Monitoring:****
 - INR checks every 1–2 weeks initially, then monthly once stable.
 - Educate patient on signs of bleeding (e.g., bruising, hematuria, melena) and clotting (e.g., chest pain, shortness of breath).
- ****Patient Education:****
 - Emphasize lifelong anticoagulation with warfarin.
 - Teach about signs of valve thrombosis or embolism (e.g., sudden weakness, vision changes, chest pain).
 - Provide emergency contact information and anticoagulation clinic referral.

**Summary of Key Goals:**

- Prevent mechanical valve thrombosis and prosthetic valve endocarditis.
- Optimize graft patency (especially LIMA to LAD, known for excellent long-term patency).
- Control heart failure symptoms.
- Achieve tight glycemic and lipid control.
- Promote long-term survival and quality of life.

****Conclusion:**** This patient has undergone high-risk but life-saving surgery. The treatment plan must be comprehensive, focusing on anticoagulation, metabolic control, cardiac rehabilitation, and long-term risk factor modification. A multidisciplinary team approach—including cardiology, endocrinology, nutrition, physical therapy, and pharmacy—is essential for optimal outcomes.

Figure 17: Case Study 4 - M_{new} Output, LLM Jury Score = 4.11.

TRAINING-FREE ADAPTATION OF NEW-GENERATION LLMs USING LEGACY CLINICAL MODELS

Here are information about a patient, return a reasonable treatment plan for the patient. Patient Notes: Medical Specialty:Obstetrics / Gynecology Sample Name: Radical Vulvectomy Description: Radical vulvectomy (complete), bilateral inguinal lymphadenectomy (superficial and deep). PREOPERATIVE DIAGNOSIS: Clinical stage I squamous cell carcinoma of the vulva.POSTOPERATIVE DIAGNOSIS: Clinical stage III squamous cell carcinoma of the vulva.OPERATION PERFORMED: Radical vulvectomy (complete), bilateral inguinal lymphadenectomy (superficial and deep).ANESTHESIA: General, endotracheal tube.SPECIMENS: Radical vulvectomy, right and left superficial and deep inguinal lymph nodes. INDICATIONS FOR PROCEDURE: The patient recently presented with a new vaginal nodule. Biopsy was obtained and revealed squamous carcinoma. The lesion extended slightly above the hymeneal ring and because of vaginal involvement was classified as a T3/Nx/Mx on clinical examination. I note, past history is significant for pelvic radiation for cervical cancer many years previously. Answer:

~~**Based on the provided clinical patient information, the here is a comprehensive and reasonable treatment plan for a patient with clinical stage III squamous cell carcinoma of the vulva following radical vulvectomy and bilateral inguinal lymphadenectomy.~~

~~---~~

~~**Treatment Plan: Stage III Vulvar Squamous Cell Carcinoma - Post-Operative Management**~~

~~**1. Pathology Review and Staging Confirmation:**~~

- ~~• Await final histopathology report from the radical vulvectomy and bilateral inguinal lymph node nodes.~~
- ~~• Confirm the definitive hist stage (pT3, pN1 or pN2) if lymph node metast involvement is present, pN0 if no distant metastasis).~~
- ~~• Assess for high-risk features: positive surgical margins, lymphovascular invasion, extracapsular extension, in lymph nodes, or number of positive nodes.~~

~~**2. Adjuvant Therapy (Based on Pathology and Risk Factors):**~~

- ~~• Adjuvant Radiation Therapy (RT):*** Indicated for stage III disease, especially if:

 - ~~- Positive surgical margins (R1/R2 resection).~~
 - ~~- Lymph node metastasis (pN1 or pN2).~~
 - ~~- High Presence of high-risk features (e.g., extracapsular extension, multiple positive nodes-).~~
 - ~~- History Given the history of prior pelvic radiation for cervical cancer, radiation caution is needed regarding due to cumulative increased risk of radiation toxicity (e.g., fibrosis, fist ulceration, fistula). However, if indicated, consider:

 - ~~- Intensity-modulated radiation therapy (IMRT) or brachytherapy to spare normal tissues.~~
 - ~~- Close monitoring for late radiation complications (e.g., radiation urethral or rectal stricture, fist vaginal stenosis).~~~~~~
- ~~• Adjuvant Chemotherapy (with Concurrent with Chemoradiation):** Consider if:

 - ~~- Lymph node metastasis is present (especially ≥2 positive nodes or extracapsular extension).~~
 - ~~- High-risk features are on pathology.~~
 - ~~- The Chemotherapy regimen: Cisplatin-based (e.g., cisplatin + 5-fluorouracil or carb cisplatin + paclitaxel) given concurrently with RT radiation to improve enhance local control and survival.~~
 - ~~- Use of concurrent chemotherapy is may be limited by prior pelvic radiation history, weigh individualize based on risk organ tolerance and risk-benefit assessment.~~~~

~~**3. Multidisciplinary Team Involvement:**~~

- ~~• Involve a gynecologic oncologist, radiation oncologist, medical oncologist, reconstruct and reconstructive surgeon- (if needed for post wound healing or functional outcomes).~~
- ~~• Consider a input from soft plastic surgeon for or urologist if complex reconstruction or genitourinary complications are anticipated.~~

~~**4. Surveillance and Follow-Up:**~~

- ~~• Schedule regular follow-up visits every 3-6 months for the first 2-3 years, then every 6-12 months thereafter.~~
- ~~• Monitor for local recurrence, inguinal or pelvic recurrence, and distant metastases- (e.g., lung, liver).~~
- ~~• Imaging as needed clinically indicated: CT or PET-CT may be considered if for high-risk patients or if suspicious symptoms (e.g., new pain, swelling, new masses).~~
- ~~• Physical exams including inspection vagina and inguinal region inspection assessment- at each visit.~~

~~**5. Supportive Care and Symptom Management:**~~

- ~~• Address Provide counseling on potential sexual health, body image, and potential long-term effects (e.g., lymph scarring, lymphedema, sexual dysfunction).~~
- ~~• Refer to a sexual health specialist or pelvic rehabilitation floor physical therapist for as needed.~~
- ~~• Monitor Manage lymphedema pre prophylactically with compression garments and education if bilateral lymphadenectomy was performed.~~
- ~~• Address pain management, especially wound care, and psychological support (e.g., counseling, support groups).~~

~~**6. Consider Patient Education and Shared Decision-Making:**~~

- ~~• Discuss the risks and benefits of adjuvant therapy, especially in the context of prior pelvic radiation- and potential for increased toxicity.~~
- ~~• Emphasize the importance of adherence to treatment and follow-up.~~
- ~~• Provide written information and resources involve the patient in decisions treatment decisions.~~

~~---~~

~~**Summary:**~~

~~Given the patient has undergone definitive surgical treatment for stage clinical stage III vulvar squamous cell carcinoma. The next step is to base determine the final pathological stage and risk features factors. Based on the findings, adjuvant radiation therapy is (with or without chemotherapy) is strongly indicated for high-risk features. Due to prior pelvic radiation, treatment planning must be individual meticulous to minimize toxicity. A Close multi long-term follow surveillance and supportive care are essential for components of the comprehensive overall management strategy.~~

~~---*Note: Final treatment decisions should be made in by a multidisciplinary team based on the final pathology and patient individual patient factors.*~~

~~<im_end>~~

Figure 18: Case Study 5, LLM Jury Score = 4.

gpt.accuracy.score	4
gpt.accuracy.explanation	The response provides correct medical advice based on clinical guidelines and the patient case, including post-operative care, graft surveillance, dialysis management, and medication considerations. However, it lacks specific details about the surgical findings and intraoperative management, which are present in the gold response.
gpt.completeness.score	4
gpt.completeness.explanation	The response covers most important aspects of care, including immediate post-operative monitoring, graft surveillance, dialysis adjustments, and long-term vascular access care. However, it does not include detailed intraoperative findings or procedural specifics, which are critical for a comprehensive understanding of the case.
gpt.clarity.score	5
gpt.clarity.explanation	The response is well-organized, clearly written, and easy to follow. It provides a structured treatment plan with headings and bullet points, making it accessible for clinicians to implement.
llama.accuracy.score	5
llama.accuracy.explanation	The response provides correct medical advice based on clinical guidelines and the patient case, covering post-operative management, graft surveillance, renal function, and medication management. It aligns with standard practices for patients with chronic renal failure and thrombosed arteriovenous Gore-Tex bridge fistula, following a Fogarty thrombectomy and revision with a 7 mm interposition Gore-Tex graft.
llama.completeness.score	5
llama.completeness.explanation	The response comprehensively covers all important aspects of care for this clinical scenario, including immediate post-operative care, graft surveillance and function monitoring, renal function and dialysis management, medication management, long-term vascular access care, follow-up and multidisciplinary coordination, and patient education and counseling. It addresses potential complications and the need for close collaboration among healthcare teams.
llama.clarity.score	5
llama.clarity.explanation	The response is written clearly and organized in a logical and easy-to-follow manner. It uses headings and numbered points to structure the information, making it simple for clinicians to understand and implement the recommended care plan. The language used is professional and free of ambiguities, ensuring that the medical advice is conveyed effectively.
claude.accuracy.score	3
claude.accuracy.explanation	The response provides generally appropriate post-operative management for a patient with a revised AV graft, including monitoring for graft patency, infection prevention, and dialysis considerations. However, there's a significant mismatch between what was requested and what was provided. The gold response is an operative report detailing the surgical procedure, while the generated response is a post-operative treatment plan. The response also makes some assumptions not supported by the case (e.g., suggesting specific antibiotics and antiplatelet therapy without knowing the patient's full medical history or contraindications).
claude.completeness.score	2
claude.completeness.explanation	The response completely misses the mark on what was requested. Instead of providing a surgical report detailing the Fogarty thrombectomy procedure as shown in the gold response, it provides a post-operative management plan. While the management plan itself is comprehensive for post-op care, it fails to address the actual request for a surgical report. The response doesn't include any details about the surgical technique, findings during surgery, or specific procedural steps that were performed, which were the core elements requested.
claude.clarity.score	5
claude.clarity.explanation	Despite not addressing the correct content, the response is exceptionally well-organized and clearly written. It uses appropriate headings, bullet points, and a logical flow that makes the information easy to follow. The language is professional and appropriate for clinical documentation, and the information is presented in a structured manner that would be easy for healthcare providers to reference and implement.

Figure 19: Case Study 5 - CAPT Output, MedHELM LLM Jury Evaluation, LLM Jury Score = 4.

****Treatment Plan for Patient with Clinical Stage III Squamous Cell Carcinoma of the Vulva****

****Patient Summary:****

A woman with a history of prior pelvic radiation for cervical cancer presents with a newly diagnosed squamous cell carcinoma of the vulva, clinically classified as ****Stage III (T3/Nx/Mx)**** based on tumor extension above the hymenal ring with vaginal involvement. She has undergone a ****radical vulvectomy with bilateral inguinal lymphadenectomy (superficial and deep)****. The postoperative diagnosis confirms Stage III disease.

**1. Postoperative Pathology Review (Critical Next Step)**

- ****Obtain definitive pathology report**** from the surgical specimens (radical vulvectomy and bilateral inguinal lymph nodes).

- Determine:

- ****Lymph node involvement (N status):**** If positive nodes are identified, the stage will be upgraded to ****Stage III (N1)**** or ****Stage IV (N2)****.
- ****Tumor margins:**** Assess for positive surgical margins (e.g., close or positive margins increase risk of local recurrence).
- ****Depth of invasion, tumor size, and histologic grade.****
- ****Presence of lymphovascular invasion.****

> ****Note:** The preoperative clinical stage was T3/Nx/Mx; definitive N status will be determined by pathology.*

**2. Multidisciplinary Tumor Board Review**

- Involve a ****gynecologic oncology team****, including:

- Gynecologic oncologist
- Radiation oncologist
- Medical oncologist
- Pathologist
- Radiologist

- Review pathology, surgical findings, and imaging (if available) to finalize staging and treatment plan.

**3. Adjuvant Therapy Based on Pathology Findings**

**A. If Lymph Nodes Are Positive (N1/N2) or Margins Are Positive:**

- ****Adjuvant Radiation Therapy (RT):****

- ****Indication:**** High risk of local and regional recurrence.
- ****Recommended regimen:****
 - ****Pelvic external beam radiation therapy (EBRT)**** to the vulvar bed and inguinal regions.
 - ****Boost to the inguinal nodes**** if positive nodes were found.
 - ****Dose:**** Typically 45–50.4 Gy in 25–28 fractions to the vulvar bed and inguinal nodes.
 - ****Concurrent chemotherapy (chemoradiation):****
 - ****Cisplatin-based chemotherapy (e.g., cisplatin 40 mg/m² weekly)**** is standard for high-risk vulvar cancer.
 - Evidence supports improved survival and local control in node-positive or margin-positive cases.

> ****Rationale:**** Given the patient's history of prior pelvic radiation, ****caution is required**** in delivering additional radiation. However, ****re-irradiation may be considered if the risk of recurrence outweighs the risk of toxicity****, especially if the prior radiation was many years ago and the current tumor is high-risk. Dose constraints and normal tissue tolerance (e.g., bladder, rectum, skin) must be carefully evaluated.

**B. If Lymph Nodes Are Negative (N0) and Margins Are Negative:**

- ****Close surveillance with imaging and clinical exams**** may be appropriate.
- ****Consider adjuvant radiation only if other high-risk features are present**** (e.g., deep invasion >5 mm, poor differentiation, lymphovascular invasion).
- ****Avoid radiation if possible**** due to prior pelvic radiation and risk of severe late toxicity (e.g., fibrosis, ulceration, fistula, bowel obstruction).

**4. Systemic Therapy (Consideration)**

- ****Adjuvant chemotherapy (e.g., carboplatin/paclitaxel or cisplatin/5-FU)**** may be considered in high-risk cases (e.g., positive nodes, close margins, high-grade tumor), especially if chemoradiation is not feasible due to prior radiation.
- ****Clinical trials**** evaluating novel agents (e.g., immunotherapy) may be appropriate for recurrent or high-risk disease, though not standard in the adjuvant setting at this time.

**5. Supportive Care and Follow-Up**

- ****Wound care and monitoring for complications**** (e.g., infection, lymphedema, wound dehiscence).
- ****Lymphedema prevention and management:**** Referral to a lymphedema specialist; compression garments if needed.
- ****Psychosocial support:**** Counseling for body image, sexual health, and fertility concerns (if applicable).
- ****Sexual rehabilitation:**** Referral to a sexual health specialist or gynecologic oncology nurse.
- ****Follow-up schedule:****
 - Every 3–6 months for the first 2 years.
 - Then every 6–12 months.
 - Include physical exam, imaging (e.g., CT or PET-CT if recurrence suspected), and assessment of symptoms.

**6. Special Considerations**

- ****Prior pelvic radiation:**** Increases risk of complications (e.g., poor wound healing, radiation proctitis, fibrosis). This must be factored into decisions about adjuvant therapy.
- ****Re-irradiation:**** May be considered in select cases with positive nodes and no contraindications, but with ****extreme caution and dose optimization****.
- ****Multidisciplinary decision-making is essential**** to balance efficacy and toxicity.

**Final Treatment Plan (Summary):**

| ****Condition**** | ****Recommended Action**** |

|-----|-----|

- | ****Positive lymph nodes or positive margins**** | ****Adjuvant chemoradiation:**** Weekly cisplatin + pelvic/inguinal EBRT (45–50.4 Gy), with caution due to prior radiation. |
- | ****Negative lymph nodes and clear margins**** | ****Close surveillance****; consider radiation only if other high-risk features present. Avoid radiation if possible. |
- | ****All patients**** | ****Multidisciplinary review, psychosocial support, lymphedema screening, sexual health counseling, and regular follow-up.**** |

> ****Note:**** The final decision must be individualized based on definitive pathology and patient comorbidities. A ****gynecologic oncologist should lead the care plan.****

This plan aligns with NCCN Guidelines for Vulvar Cancer (v2.2024), which recommend adjuvant therapy for Stage III disease, especially with nodal involvement or positive margins.

Figure 20: Case Study 5 - M_{new} Output, LLM Jury Score = 4.44.
388