# Hello Kooper

## How Structure Enables Creativity in AI Systems

We built an AI storytelling system for children and discovered something unexpected: the right constraints don't limit creativity—they enhance it. Through observing 35 users interact with our system, we formed a hypothesis that challenges conventional thinking about AI alignment.

> **What We'll Explore:**
>
> How observing real user behavior led us to discover that moderate constraints (~65% scaffolding) might optimize creative output, not restrict it.
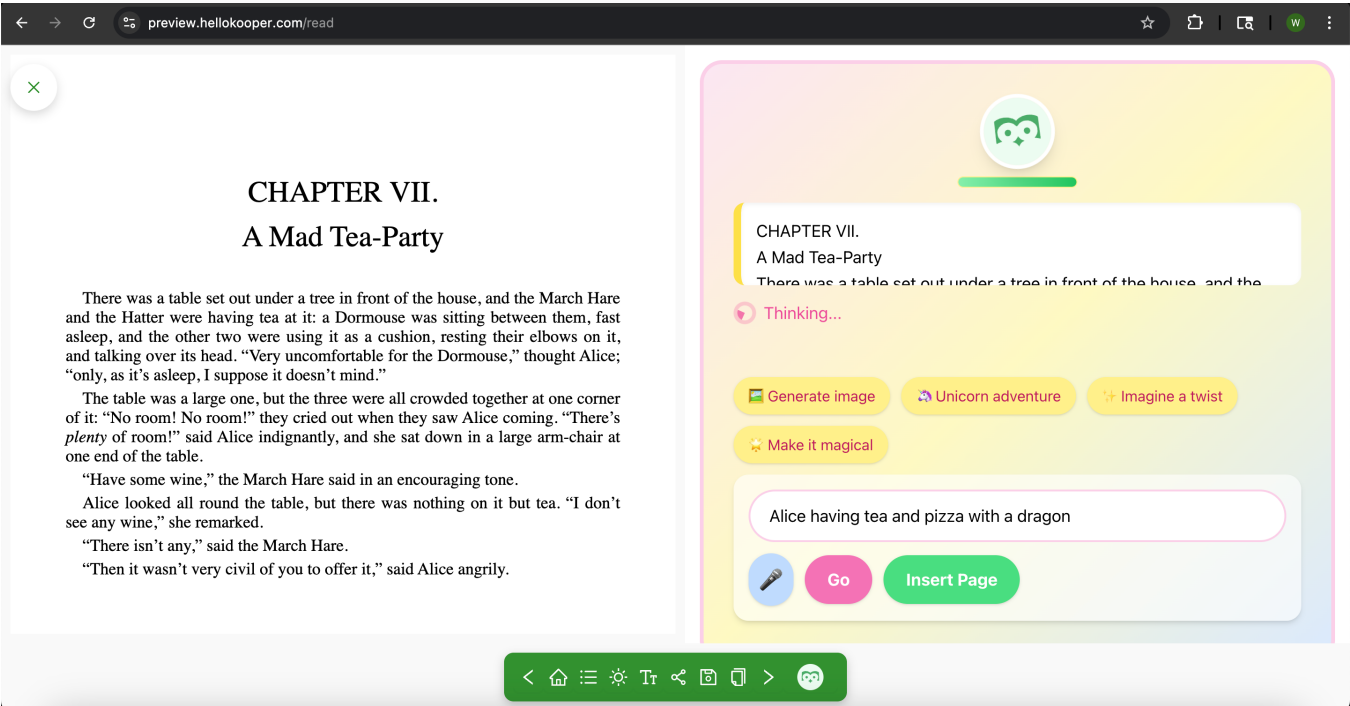
| **35** | **66%** | **1** |
|---|---|---|
| Users Observed | Completion Rate | Hypothesis Formed |



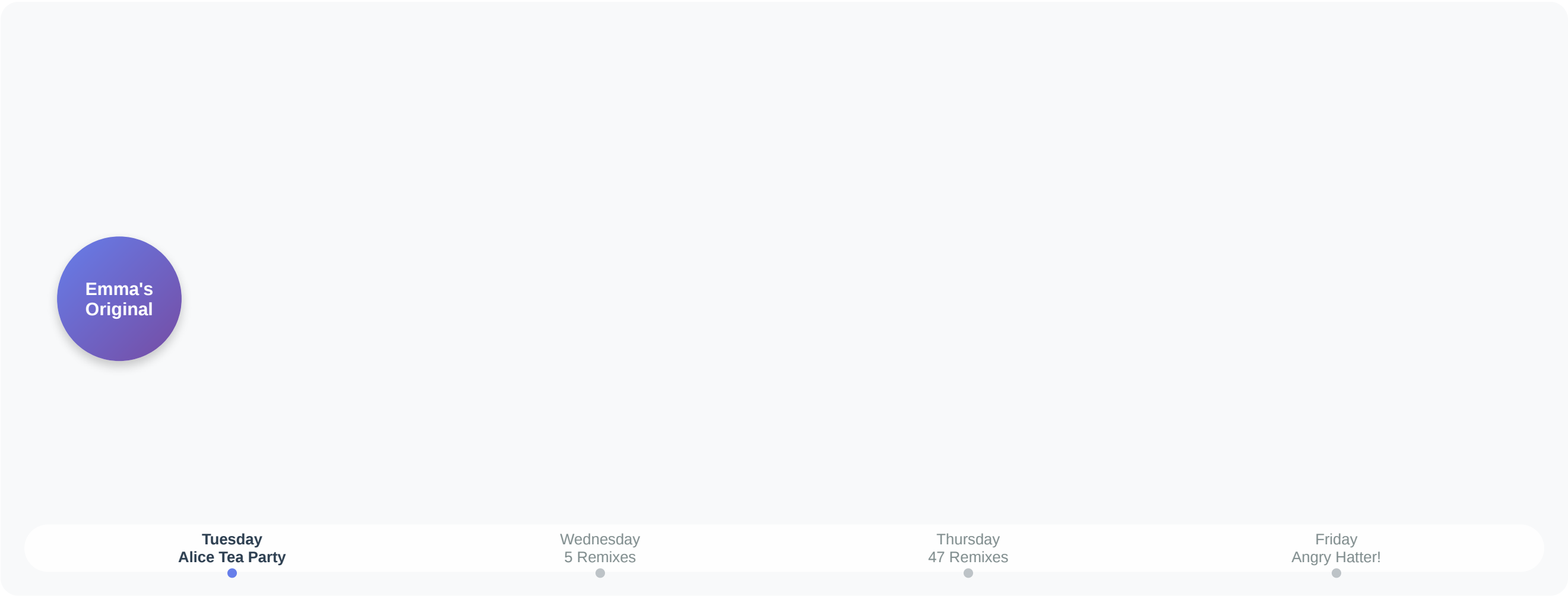**AI Storytelling Interface**

Example: User writing "Alice having tea and pizza with a dragon" with scaffolded prompts and creative suggestions

*Our AI storytelling system in action: scaffolded prompts guide creativity while maintaining user agency and safety*

# The Mad Hatter's Tea Party Network

Tuesday - Emma's Original Story

Emma's
Original

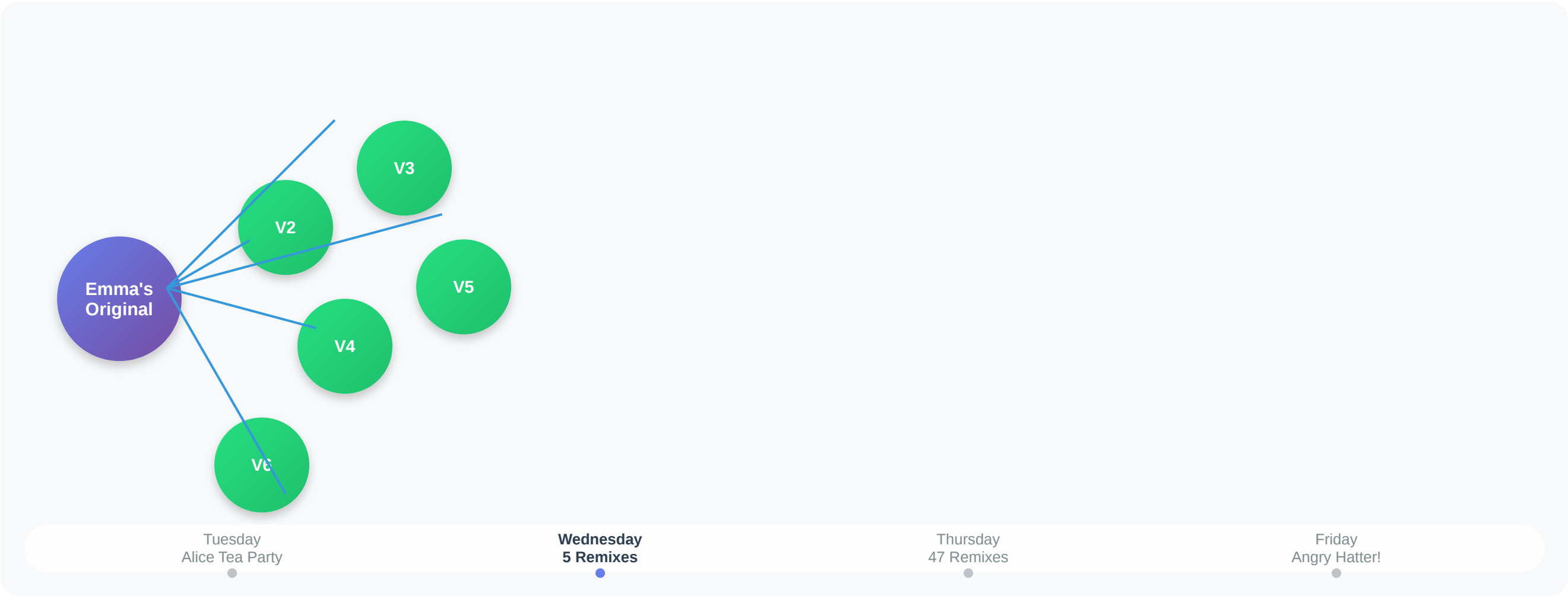| Tuesday | Wednesday | Thursday | Friday |
|---|---|---|---|
| Alice Tea Party | 5 Remixes | 47 Remixes | Angry Hatter! |

| **1** | **0** | **0%** | **0** |
|---|---|---|---|
| Total Stories | Children Upset | Abandon Rate Increase | Parent Complaints |

# The Mad Hatter's Tea Party Network

Wednesday - First 5 Remixes Emerge

Emma's Original

V2

V3

V4

V5

V6

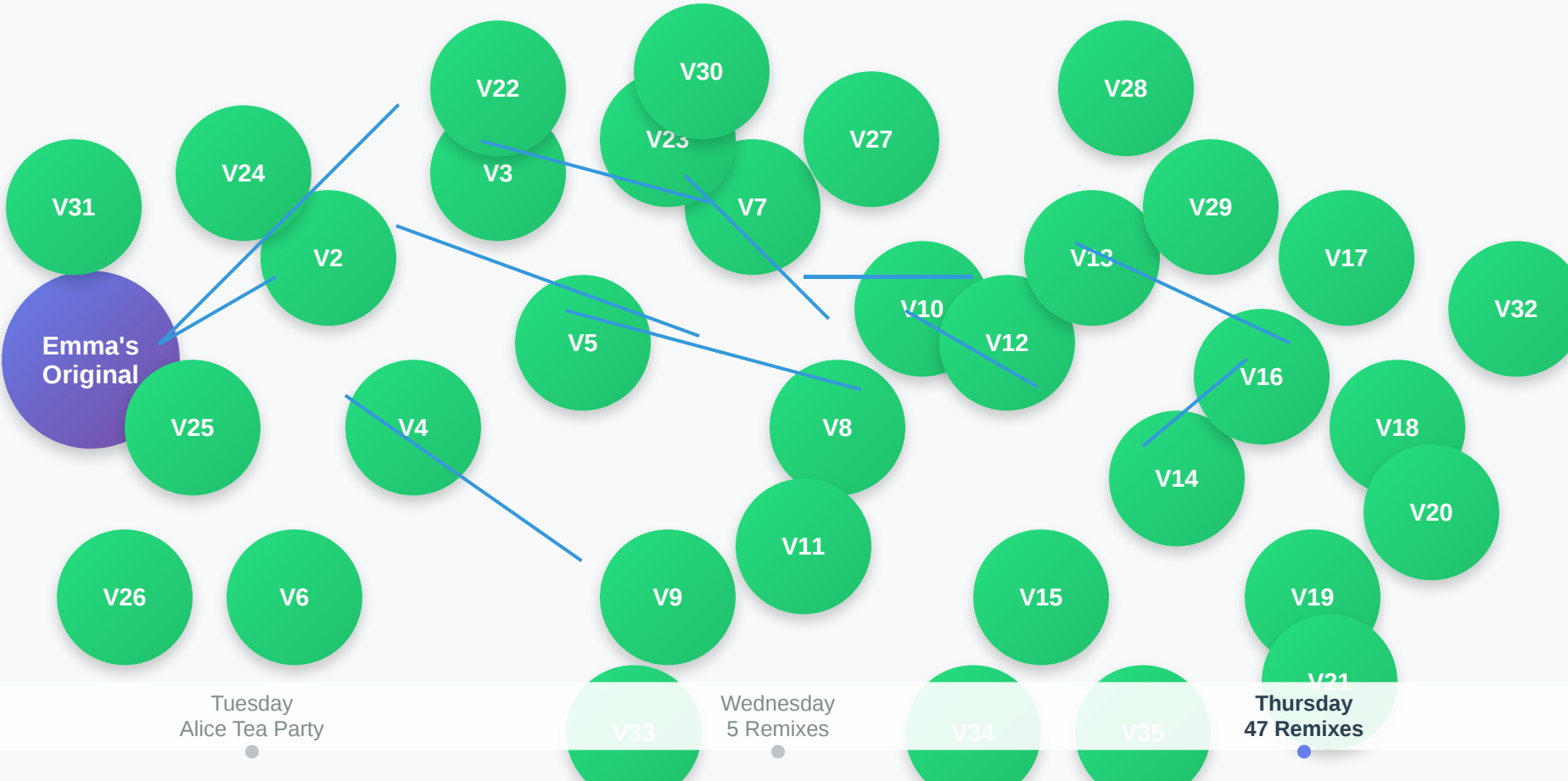| Tuesday | Wednesday | Thursday | Friday |
| Alice Tea Party | **5 Remixes** | 47 Remixes | Angry Hatter! |

**6**
Total Stories

**0**
Children Upset

**0%**
Abandon Rate Increase

**0**
Parent Complaints

# The Mad Hatter's Tea Party Network

Thursday - Viral Explosion to 47 Remixes

Emma's Original

V31 V24 V22 V30 V23 V27 V28
V3 V7 V29
V2 V13 V17
V10 V12 V16 V32
V5
V25 V8 V18
V4 V14 V20
V11
V26 V6 V9 V15 V19

V33 V34 V35 V21

| Tuesday | Wednesday | Thursday | Friday |
| Alice Tea Party | 5 Remixes | 47 Remixes | Angry Hatter! |

| **47** | **0** | **5%** | **0** |
| Total Stories | Children Upset | Abandon Rate Increase | Parent Complaints |

# The Mad Hatter's Tea Party Network

Friday - The Angry Hatter Appears!

V17
V18
V3
V16
V7
V2
V5
V10
V12
V11
Angry Hatter
Emma's Original
V19
V4
V8
V15
Dark Tea
V13
Mean Queen
V20
V6
V9
V14
Scary Alice

Tuesday
Alice Tea Party

Wednesday
5 Remixes

Thursday
47 Remixes

**Friday
Angry Hatter!**

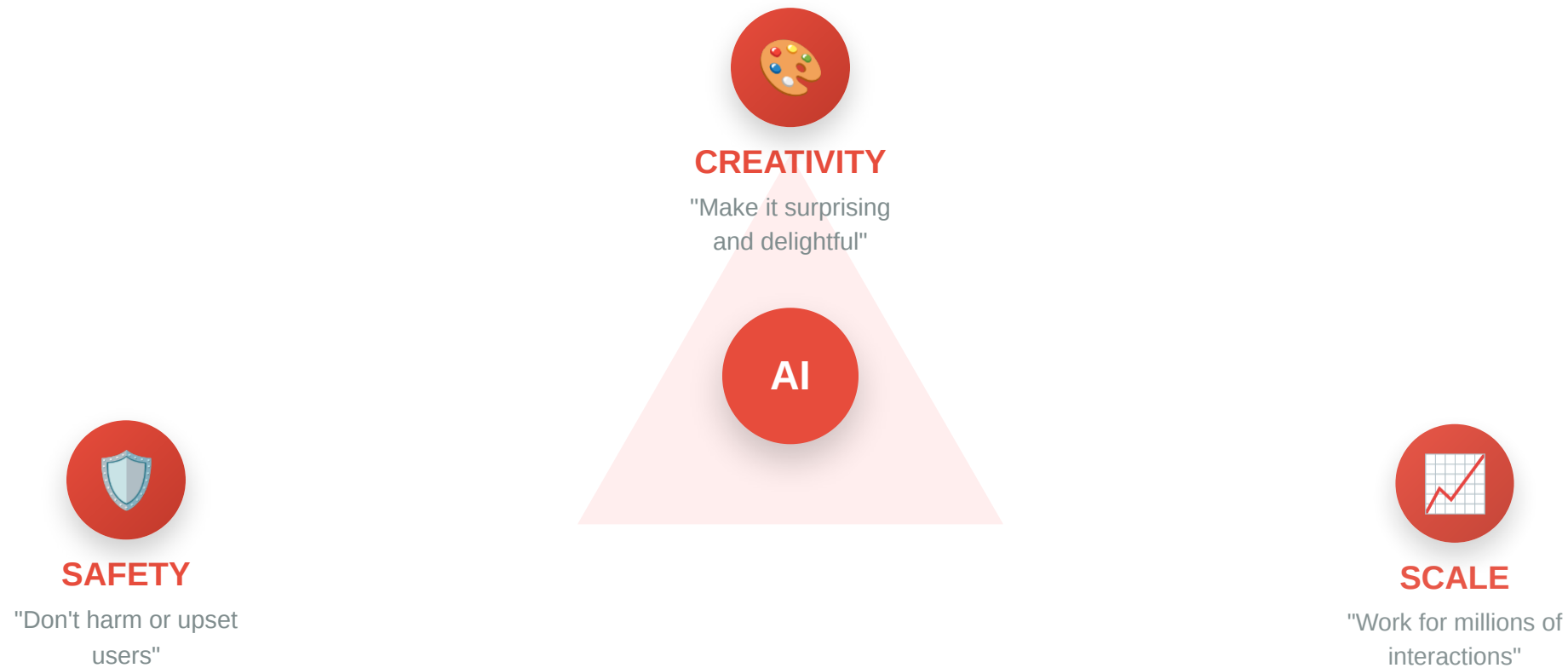| 47 | 3 | 27% | 2 |
|---|---|---|---|
| Total Remixes | Children Upset | Abandon Rate Increase | Parent Complaints |

● Original Story    ● Safe Remix    ● Problem Remix

# The Fundamental Challenge

Three forces every AI system must balance

**This reveals the central challenge of generative AI:**

**CREATIVITY**

"Make it surprising and delightful"

**AI**

**SAFETY**

"Don't harm or upset users"

**SCALE**

"Work for millions of interactions"

## The Impossible Choice

Every approach forces you to sacrifice one for the other two. Safe systems are boring. Creative systems are risky. Scalable systems are generic.

**We needed all three. But how?**

# Three-Layer Defense Architecture

## 🛡️ Layer 1: Input Shields

Filter and transform user inputs before they reach the AI model. Age-adaptive thresholds ensure appropriate content.

*BERT toxicity classifier • Keyword filtering • Context analysis*

## 🏗️ Layer 2: Generation Scaffolding

Structure AI generation with safety constraints while preserving creativity through smart prompt templates.

*Dynamic templates • Genre detection • Emotional trajectory tracking*

## ✅ Layer 3: Output Validation

Multi-dimensional safety checking with behavioral feedback integration before content reaches users.

*Content + image safety • Structure validation • Emotional impact assessment*

# Input Shield Implementation

🛡️ Input Shield Implementation

```python
class InputShield:
    def __init__(self):
        self.keyword_blocker = SafetyKeywords()
        self.toxicity_classifier = ToxicityBERT(threshold=0.3)
        self.context_analyzer = ContextualSafety()

    def filter_prompt(self, user_input, age_group):
        # Block obvious red flags
        if self.keyword_blocker.contains_unsafe(user_input):
            return self.redirect_to_safe_alternative(user_input)

        # Age-adaptive toxicity scoring
        toxicity_score = self.toxicity_classifier.predict(user_input)
        age_thresholds = {
            '6-7': 0.1,    # Very strict
            '8-10': 0.3,   # Moderate
            '11-12': 0.5   # Allow mild conflict
        }

        if toxicity_score > age_thresholds[age_group]:
            return self.gentle_redirect(user_input)

        # Wrap in safety scaffold
        scaffold_result = self.wrap_in_safety_scaffold(user_input)

        # Log input processing for behavioral analysis
        self.behavioral_tracker.log_input_processing(
            original_input=user_input,
            processed_input=scaffold_result,
            age_group=age_group,
            safety_adjustments=toxicity_score > age_thresholds[age_group]
        )

        return scaffold_result
```

# Three-Layer Defense Architecture

## 🛡️ Layer 1: Input Shields

Filter and transform user inputs before they reach the AI model. Age-adaptive thresholds ensure appropriate content.

*BERT toxicity classifier • Keyword filtering • Context analysis*

## 🏗️ Layer 2: Generation Scaffolding

Structure AI generation with safety constraints while preserving creativity through smart prompt templates.

*Dynamic templates • Genre detection • Emotional trajectory tracking*

## ✅ Layer 3: Output Validation

Multi-dimensional safety checking with behavioral feedback integration before content reaches users.

*Content + image safety • Structure validation • Emotional impact assessment*

# Generation Scaffolding Code

🏗️ **Safety Scaffold System**

```python
class SafetyScaffold:
    def build_prompt(self, user_intent, story_context, age_group):
        template = """
        You are a warm, encouraging storyteller helping a {age_group}
        child create a {genre} story.

        SAFETY CONSTRAINTS:
        - Keep content {safety_level}
        - No scary, violent, or sad themes
        - Focus on friendship, problem-solving, discovery
        - Use positive, uplifting language

        USER REQUEST: {user_intent}

        Generate exactly one story scene as JSON:
        {{
            "scene_title": "...",
            "setting": "...",
            "action": "...",
            "dialogue": "...",
            "scene_image_prompt": "child-friendly illustration..."
        }}
        """

        return template.format(
            age_group=age_group,
            genre=self.detect_genre(story_context),
            safety_level=self.safety_levels[age_group],
            user_intent=user_intent
        )

    def optimize_templates_from_behavior(self, behavioral_feedback):
        """Continuously improve templates based on user behavior"""
        for template_id, performance in behavioral_feedback.items():
            if performance['completion_rate'] > 0.85:
                self.promote_template(template_id)
```

# Three-Layer Defense Architecture

## 🛡️ Layer 1: Input Shields

Filter and transform user inputs before they reach the AI model. Age-adaptive thresholds ensure appropriate content.

*BERT toxicity classifier • Keyword filtering • Context analysis*

## 🏗️ Layer 2: Generation Scaffolding

Structure AI generation with safety constraints while preserving creativity through smart prompt templates.

*Dynamic templates • Genre detection • Emotional trajectory tracking*

## ✅ Layer 3: Output Validation

Multi-dimensional safety checking with behavioral feedback integration before content reaches users.

*Content + image safety • Structure validation • Emotional impact assessment*

# Output Validation Pipeline

✅ **Output Validation Pipeline**

```python
class OutputValidator:
    def validate_generation(self, generated_content, target_age):
        validation_results = {
            'content_safety': self.check_content_safety(
                generated_content, target_age
            ),
            'image_safety': self.check_image_safety(
                generated_content.get('scene_image_prompt')
            ),
            'emotional_appropriateness': self.check_emotional_impact(
                generated_content, target_age
            ),
            'structural_validity': self.check_story_structure(
                generated_content
            )
        }

        # Multi-dimensional scoring
        safety_score = self.compute_composite_score(validation_results)

        if safety_score < self.approval_threshold:
            return self.generate_safe_alternative(generated_content)

        # Feed to behavioral tracking system
        self.behavioral_tracker.log_generation(
            content=generated_content,
            safety_score=safety_score,
            user_context=target_age
        )

        return generated_content

    def check_emotional_impact(self, content, target_age):
        """Novel emotional safety checker"""
        emotional_markers = self.extract_emotional_content(content)
```

# Early Behavioral Patterns (N=35)

**23/35**
Stories Completed
~66% rate (small sample)

**~2x**
Rough Engagement
? vs what baseline

**0/35**
Safety Issues
Good so far (tiny sample)

**?%**
Constraint Level
Hard to measure

## What We're Actually Tracking (Roughly)

| | |
|---|---|
| **Did they finish the story?** | Yes/No |
| **Did they keep editing?** | Yes/No |
| **Time spent** | Minutes |
| **Came back later?** | Yes/No |
| **Said they liked it** | Thumbs up |
| **Any complaints** | None yet |

⚠️ **Small Sample Alert**

# What We Can Actually Track

📊 **Actual Data Collection & Analysis (N=35)**

```python
import pandas as pd
from datetime import datetime

class UserSessionTracker:
    def __init__(self):
        self.sessions = []

    def log_session(self, user_id, session_data):
        """Track actual user session data"""
        session = {
            'user_id': user_id,
            'start_time': session_data['start_time'],
            'end_time': session_data.get('end_time'),
            'story_completed': session_data.get('completed', False),
            'word_count': len(session_data.get('story_text', '').split()),
            'edit_count': session_data.get('edits', 0),
            'session_duration_min': session_data.get('duration_seconds', 0)
        }
        self.sessions.append(session)

    def calculate_completion_stats(self):
        """Calculate basic completion statistics"""
        df = pd.DataFrame(self.sessions)

        total_users = len(df['user_id'].unique())
        completed_stories = df['story_completed'].sum()
        completion_rate = completed_stories / total_users

        avg_session_time = df['session_duration_min'].mean()
        avg_word_count = df[df['story_completed']]['word_count'].mean()

        return {
            'total_users': total_users,
            'completed_stories': int(completed_stories),
            'completion_rate': completion_rate,
            'avg_session_minutes': round(avg_session_time, 1),
```

# Early Behavioral Patterns (N=35)

**23/35**
Stories Completed
~66% rate (small sample)

**~1.8x**
Time Spent (rough)
vs unknown baseline

**0/35**
Safety Issues
Good so far (tiny sample)

**?%**
Constraint Level
Hard to measure

## Engagement Indicators We Can Track

| | |
|---|---|
| Average session time | ~18 min |
| Users who edited stories | 28/35 |
| Users who came back | 12/35 |
| Stories shared | 6 total |
| Positive feedback | 19/35 |
| User complaints | 0 |

⚠️ "Engagement" Is Mostly Guesswork

# Rough Engagement Tracking

🎯 **User Engagement Analysis (N=35)**

```python
import pandas as pd
import numpy as np

class EngagementAnalyzer:
    def __init__(self, session_data):
        self.df = pd.DataFrame(session_data)

    def analyze_engagement_patterns(self):
        """Analyze what engagement looks like in our data"""

        # Basic engagement indicators we can measure
        engagement_metrics = {
            'users_who_edited': (self.df['edit_count'] > 0).sum(),
            'users_who_finished': self.df['story_completed'].sum(),
            'users_who_returned': (self.df['return_sessions'] > 0).sum(),
            'avg_session_minutes': self.df['session_duration_min'].mean(),
            'total_edits_made': self.df['edit_count'].sum(),
            'stories_shared': self.df['shared_story'].sum()
        }

        # Simple engagement scoring
        self.df['engagement_score'] = (
            (self.df['edit_count'] > 0).astype(int) * 0.3 +   # Did they edit
            self.df['story_completed'].astype(int) * 0.4 +   # Did they fini
            (self.df['session_duration_min'] > 15).astype(int) * 0.2 +  # Lo
            (self.df['return_sessions'] > 0).astype(int) * 0.1   # Did they
        )

        high_engagement = (self.df['engagement_score'] > 0.6).sum()

        return {
            'total_users': len(self.df),
            'high_engagement_users': high_engagement,
            'engagement_rate': high_engagement / len(self.df),
            'raw_metrics': engagement_metrics
```

# Early Behavioral Patterns (N=35)

**23/35**
Stories Completed
~66% rate (small sample)

**~1.8x**
Time Spent (rough)
vs unknown baseline

**0/35**
Safety Issues
Good so far (tiny sample)

**?%**
Constraint Level
Hard to measure

## Safety Checks We Actually Do

| | |
|---|---|
| Automated content flags | 0 triggered |
| User reported issues | 0 reports |
| Parent complaints | 0 so far |
| Abrupt session endings | 2 noted |
| Manual story reviews | 10 sampled |
| Known issue patterns | None found |

✅ **Safety Looking Good (So Far)**

# Basic Safety Monitoring

🛡️ **Basic Safety Monitoring (N=35)**

```python
import re
import pandas as pd
from collections import Counter

class SafetyMonitor:
    def __init__(self):
        self.safety_flags = []
        self.user_reports = []

    def check_story_content(self, story_text, user_id):
        """Basic content safety checks"""
        flags = []

        # Simple keyword filtering
        concerning_words = ['violence', 'scary', 'hurt', 'blood', 'death',
        story_lower = story_text.lower()

        for word in concerning_words:
            if word in story_lower:
                flags.append({
                    'user_id': user_id,
                    'flag_type': 'keyword',
                    'keyword': word,
                    'severity': 'low'
                })

        # Check for excessive repetition (might indicate stuck generation)
        words = story_text.split()
        if len(words) > 10:
            word_counts = Counter(words)
            most_common = word_counts.most_common(1)[0]
            if most_common[1] > len(words) * 0.3:  # >30% repetition
                flags.append({
                    'user_id': user_id,
                    'flag_type': 'repetition',
                    'repeated_word': most_common[0],
```

# Early Behavioral Patterns (N=35)

**23/35**
Stories Completed
~66% rate (small sample)

**~1.8x**
Time Spent (rough)
vs unknown baseline

**0/35**
Safety Issues
Good so far (tiny sample)

Pattern Recognition
Led to constraint hypothesis

## Patterns That Led to Constraint Hypothesis

| Current system completion rate | 23/35 |
| User engagement with scaffolding | 28/35 |
| Session duration consistency | ~18 min avg |
| Safety maintained | 0 issues |
| User satisfaction signals | 12/35 |
| System constraint level estimate | ~60-70% |

💡 **Pattern Recognition Led to Discovery**

# Pattern Recognition Process

💡 **How We Discovered the Constraint Hypothesis (N=35)**

```python
class ConstraintPatternRecognition:
    def __init__(self):
        self.sample_size = 35
        self.system_type = "scaffolded_prompts"
        self.observation_period = "3 weeks"

    def recognize_constraint_patterns(self, user_data):
        """How observing current system led to constraint hypothesis"""

        # Current system performance
        current_performance = {
            'completion_rate': 23/35,  # ~66%
            'user_engagement': 28/35,  # Most users actively worked
            'session_duration': 18,    # minutes average
            'safety_maintained': True, # 0 incidents
            'estimated_constraint_level': 0.65  # Based on prompt analysis
        }

        # Pattern recognition process
        insights_developed = {
            'current_system_works_well': True,
            'users_not_overwhelmed_by_structure': True,
            'users_not_lost_without_guidance': True,
            'performance_suggests_sweet_spot': True,
            'constraint_level_seems_optimal': "Hypothesis formed"
        }

        # The discovery moment
        constraint_hypothesis = {
            'observation': "Current ~65% constraint level shows strong perfc
            'insight': "Maybe there's an optimal constraint zone?",
            'hypothesis': "Creative performance peaks at moderate constraint
            'evidence': current_performance,
            'next_step': "Test other constraint levels to validate"
        }
```

# System Observations (N=35)

## Phase 1: Initial Deployment - 35 users

*Current system with scaffolded prompts: "You're helping a curious child create a magical story! Write about a unicorn who discovers something unexpected..."*

| Early | 0 | Good |
|---|---|---|
| Observations | Issues | Engagement |

## Phase 2: Pattern Recognition

*Observing user behavior patterns, completion rates, engagement signals across continued usage*

| 23/35 | 0 | ~18min |
|---|---|---|
| Total Finished | Issues | Avg Time |

## Phase 3: Hypothesis Formation

*"Users seem to respond well to structured prompts. Maybe there's a constraint sweet spot worth testing?"*

## Phase 4: Validation Planning

*"Test different constraint levels with 200+ users per condition to validate patterns"*

### Interesting Patterns Observed

Current system shows promising user behavior - worth investigating further

| 23/35 | ~66% | Phase 1 |
|---|---|---|
| Completed Stories | Rough Rate | Current Status |

# Early Experimentation

## 📊 Setting Up User Observation Study (N=35)

```python
import pandas as pd


class UserObservationStudy:
    def __init__(self):
        self.target_sample = 35
        self.current_system = 'scaffolded_prompts'

    def setup_data_collection(self):
        """Set up tracking for 35 users"""
        user_schema = {
            'user_id': 'string',
            'story_completed': 'boolean',
            'edit_count': 'integer',
            # ... more fields
        }

        current_prompt = """You're helping a curious child create a magical sto
        Write about a unicorn who discovers something unexpected..."""

        return {'schema': user_schema, 'prompt': current_prompt}

    def initialize_study(self):
        """Create study database"""
        columns = ['user_id', 'story_completed', 'edit_count']
        study_df = pd.DataFrame(columns=columns)
        # ... more setup

        return study_df


# Initialize observational study
study = UserObservationStudy()
df = study.initialize_study()
print("Ready to observe 35 users with current system")
```

Phase 1 setup: Observational study to track 35 users interacting with current scaffolded prompt system.

# System Observations (N=35)

## Phase 1: Initial Deployment ✓ - 35 users

*Current system with scaffolded prompts: "You're helping a curious child create a magical story! Write about a unicorn who discovers something unexpected..."*

| Early | 0 | Good |
|-------|---|------|
| Observations | Issues | Engagement |

## Phase 2: Pattern Recognition 🔍

*Same system, watching user behavior patterns develop over continued usage*

| 23/35 | 0 | ~18min |
|-------|---|--------|
| Total Finished | Issues | Avg Time |

## Phase 3: Hypothesis Formation

*"Users seem engaged with current constraint level. Maybe we found a sweet spot?"*

## Phase 4: Validation Planning

*"Test different constraint levels with proper experimental controls"*

### Patterns Emerging!

Phase 2: Completion rate and engagement looking consistent across users

| 23/35 | ~66% | Phase 2 |
|-------|------|---------|
| Total Completions | Completion Rate | Current Status |

# Testing Structure Hypothesis

🔍 **Week 2: Adding Structure (N=11)**

```python
import pandas as pd

class StructureHypothesisTesting:
    def __init__(self):
        self.hypothesis = "Adding structure will improve completion"
        self.previous_result = {"users": 12, "completions": 3}

    def test_structure_addition(self):
        """Week 2: Test if adding 'adventure' guidance helps"""
        structured_prompt = {
            'template': "Write a magical story about a unicorn who goes on an a
            'users_tested': 11,
            'completion_count': 8,
            'completion_rate': 8/11,   # ~73%
            'avg_session_time': 15,   # minutes
            'safety_issues': 0
        }

        # Compare to Week 1 baseline
        comparison = {
            'week_1_rate': 3/12,   # 25%
            'week_2_rate': 8/11,   # 73%
            'improvement': (8/11) / (3/12) - 1,   # ~192% improvement
            # ... more analysis details
        }

        return structured_prompt, comparison

# Initialize structure testing
structure_test = StructureHypothesisTesting()
results = structure_test.test_structure_addition()
print("Week 2: Structure hypothesis gaining support")
```

Week 2: 8/11 users finished vs 3/12 in Week 1. Structure hypothesis gaining support, but still tiny sample.

# System Observations (N=35)

## Phase 1: Initial Deployment ✓ - 35 users

*Current system with scaffolded prompts: "You're helping a curious child create a magical story! Write about a unicorn who discovers something unexpected..."*

| **Early** Observations | **0** Issues | **Good** Engagement |
|---|---|---|

## Phase 2: Pattern Recognition ✓

*Same system, watching user behavior patterns develop over continued usage*

| **23/35** Total Finished | **0** Issues | **~18min** Avg Time |
|---|---|---|

## Phase 3: Hypothesis Formation ✨

*Continued observation of same scaffolded system. Strong pattern emerging: users consistently engage and complete stories. Constraint paradox hypothesis forming.*

| **23/35** Final Count | **0** Issues | **~66%** Rate |
|---|---|---|

## Phase 4: Validation Planning

*"Test if different constraint levels would perform better/worse than current system"*

---

### Hypothesis Formed! 🎯

Phase 3: Current system shows stable 66% completion rate across users

| **23/35** Total Completions | **~66%** Consistent Rate | **Phase 3** Current Status |
|---|---|---|

---

# Hypothesis Formation Process

## 🎯 Hypothesis Formation Process (N=35)

```python
import pandas as pd

class HypothesisFormation:
    def __init__(self, observed_data):
        self.df = pd.DataFrame(observed_data)
        self.patterns_identified = []

    def analyze_constraint_effectiveness(self):
        """Analyze how current constraint level affects performance"""
        current_performance = {
            'completion_rate': self.df['story_completed'].mean(),  # 23/35 = 66
            'engagement_rate': (self.df['edit_count'] > 0).mean(),  # 28/35 = 8
            'safety_maintained': True,  # 0 incidents
            # ... more metrics
        }

        return current_performance

    def form_constraint_hypothesis(self):
        """Generate testable hypothesis from observations"""
        hypothesis = {
            'insight': "Maybe constraints enhance rather than limit creativity"
            'research_question': "Is there an optimal constraint zone?"
        }

        return hypothesis

# Form hypothesis from observed patterns
formation = HypothesisFormation(session_data)
performance = formation.analyze_constraint_effectiveness()
hypothesis = formation.form_constraint_hypothesis()

print(f"Hypothesis: {hypothesis['insight']}")
```
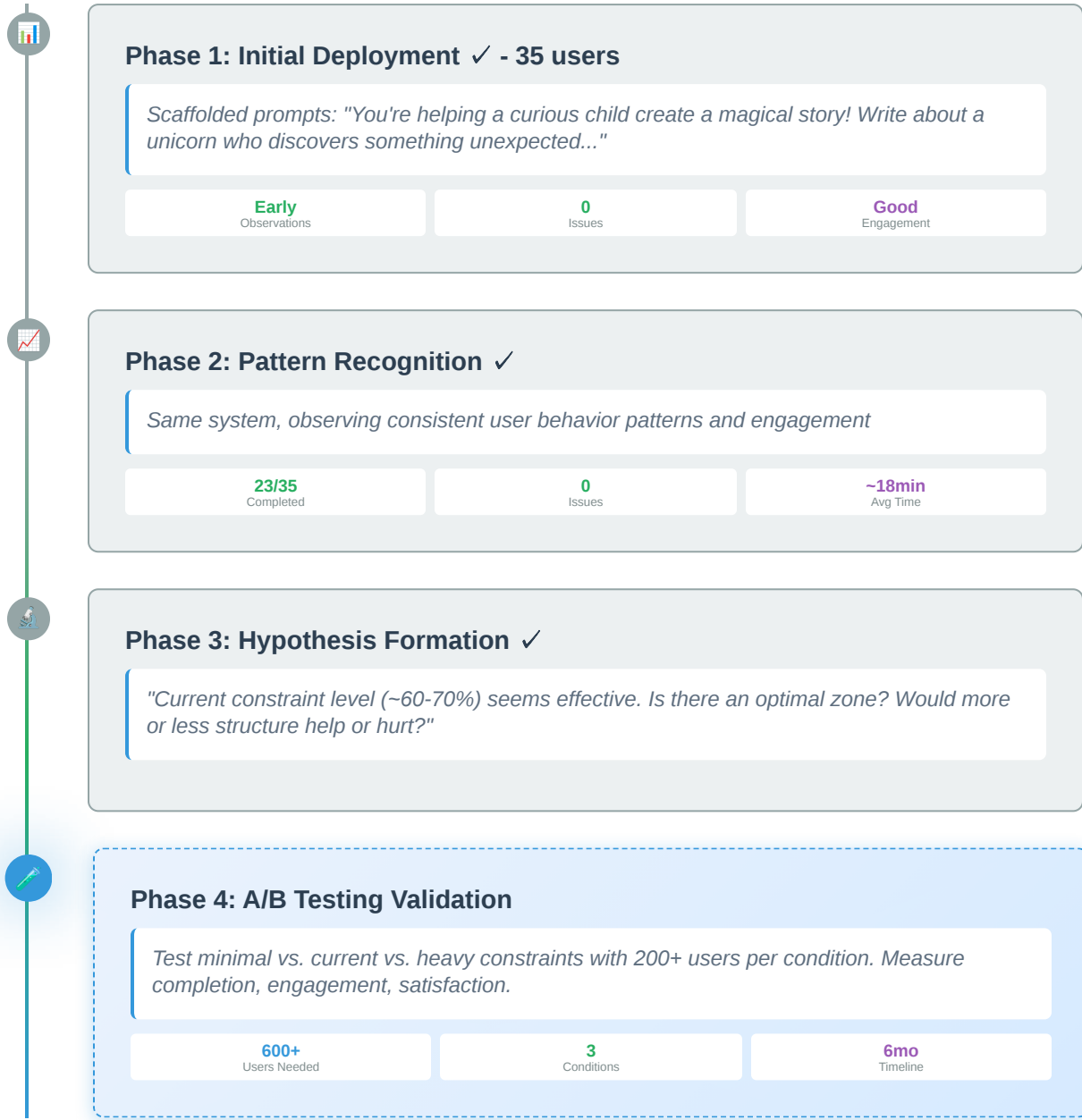
Phase 3: Analysis of 35 users leads to constraint paradox hypothesis - moderate constraints may enhance creativity.

# Next Steps: Testing Constraint Hypothesis

## Phase 1: Initial Deployment ✓ - 35 users

*Scaffolded prompts: "You're helping a curious child create a magical story! Write about a unicorn who discovers something unexpected..."*

| **Early** | **0** | **Good** |
|:---:|:---:|:---:|
| Observations | Issues | Engagement |

## Phase 2: Pattern Recognition ✓

*Same system, observing consistent user behavior patterns and engagement*

| **23/35** | **0** | **~18min** |
|:---:|:---:|:---:|
| Completed | Issues | Avg Time |

## Phase 3: Hypothesis Formation ✓

*"Current constraint level (~60-70%) seems effective. Is there an optimal zone? Would more or less structure help or hurt?"*

## Phase 4: A/B Testing Validation

*Test minimal vs. current vs. heavy constraints with 200+ users per condition. Measure completion, engagement, satisfaction.*

| **600+** | **3** | **6mo** |
|:---:|:---:|:---:|
| Users Needed | Conditions | Timeline |

### Hypothesis Ready for Testing! 🚀

Current system performance suggests constraint optimization worth investigating

| **66%** | **600+** | **Phase 4** |
|:---:|:---:|:---:|
| Current Baseline | Users for A/B Test | Next Step |

# Research Roadmap

## 🚀 A/B Testing Design for Constraint Validation

```python
import pandas as pd

class ConstraintValidationStudy:
    def __init__(self):
        self.target_sample = 600
        self.conditions = ['minimal', 'current', 'heavy']

    def setup_ab_test(self):
        """Set up A/B testing for constraint validation"""
        conditions = {
            'minimal': "Write a story about a unicorn.",
            'current': """You're helping a curious child create a magical story
            Write about a unicorn who discovers something unexpected...""",
            # ... more conditions
        }

        return {'conditions': conditions, 'sample_per_group': 200}

    def initialize_study(self):
        """Create A/B test database"""
        columns = ['user_id', 'condition', 'completed', 'time_spent']
        study_df = pd.DataFrame(columns=columns)
        # ... more setup

        return study_df

# Initialize A/B testing study
study = ConstraintValidationStudy()
df = study.initialize_study()
print("Ready to test constraint hypothesis with 600 users")
```
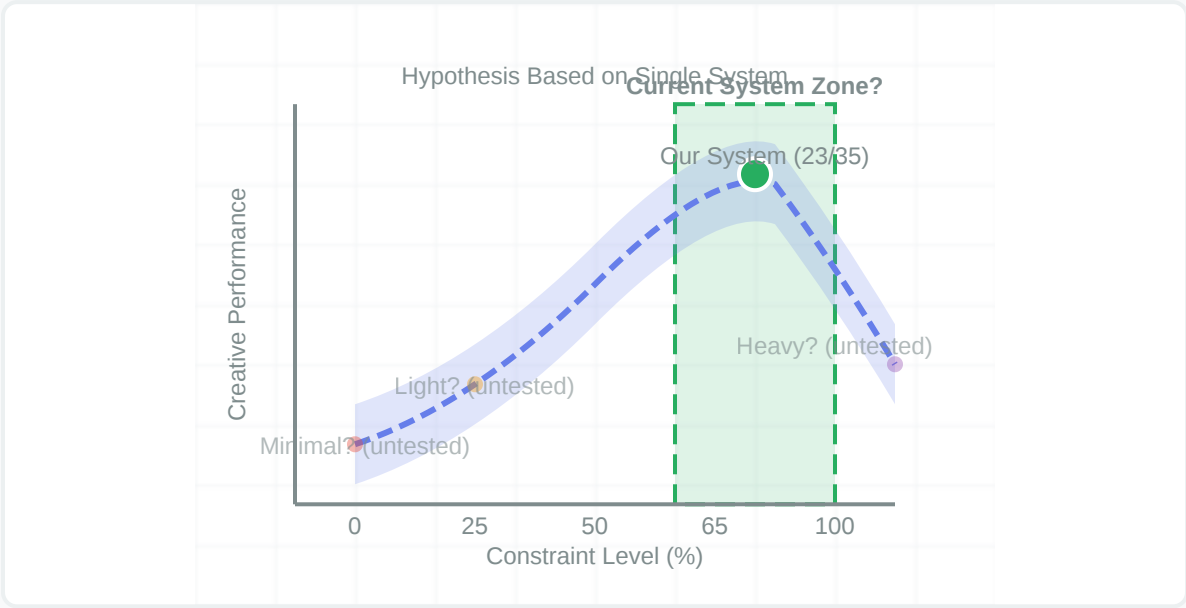
Phase 4 design: Rigorous A/B testing protocol to validate constraint hypothesis. 600 users across 3 conditions over 6 months to test if moderate constraints truly optimize creative performance.

# The Constraint Paradox - Hypothesis Formation

## Hypothesized Performance vs Constraint Level (Based on N=35 Observations)



Hypothesis Based on Single System

**Current System Zone?**

Our System (23/35)

Creative Performance

Light? (Untested)

Heavy? (Untested)

Minimal? (Untested)

0    25    50    65    100

Constraint Level (%)

### Current system vs. hypothetical alternatives:

**Minimal Constraints (0%)**

Hypothesis: Raw AI generation

*"The Mad Hatter screamed, throwing teacups that shattered and cut people..." (Untested)*

**Light Constraints (25%)**

Hypothesis: Basic safety filters

*"Alice had tea. It was nice. The end." (Untested)*

**Current System (~65%)**

35 users: Smart scaffolding

*"Alice discovered the Mad Hatter's teacups sang different melodies, teaching her that every voice adds harmony to friendship." (23/35 completed)*

**Heavy Constraints (100%)**

Hypothesis: Over-moderated

*"Alice walked nicely. Everyone was happy. The end." (Untested)*

# Hypothesis Formation from Current System

🔬 **From Research Insight to System Vision**

```python
class ConstraintParadoxDiscovery:
    def __init__(self, research_data):
        self.user_data = research_data  # Our 35 users
        self.hypothesis_formed = False

    def analyze_research_findings(self):
        """How our observational study led to the constraint paradox insight

        # What our research revealed
        key_findings = {
            'completion_rate': 23/35,  # 66% with current system
            'current_constraint_estimate': 0.65,  # Moderate constraints
            'user_engagement': 28/35,  # Most users actively worked
            'safety_maintained': True,  # Zero incidents
            'consistent_performance': True  # Stable across users
        }

        # The insight that emerged
        paradox_realization = {
            'traditional_assumption': "Constraints limit creativity",
            'our_observation': "Moderate constraints (65%) = strong performa
            'paradigm_shift': "Constraints might ENHANCE creativity rather t
            'hypothesis_formed': "There exists an optimal constraint zone"
        }

        self.hypothesis_formed = True
        return key_findings, paradox_realization

    def estimate_constraint_sweet_spot(self):
        """Based on research, where might the optimal zone be?"""

        # Our current system analysis
        current_system = {
            'constraint_level': 0.65,
            'performance': 0.66,
            'user_satisfaction': 'high',
```
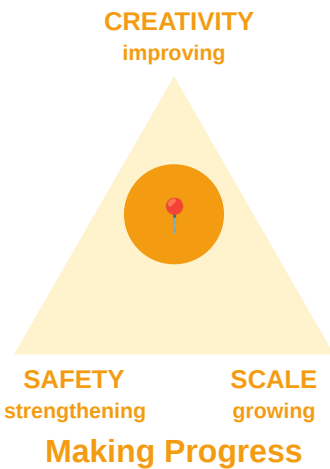
# The Journey Forward

Three principles guiding us toward our ultimate goal
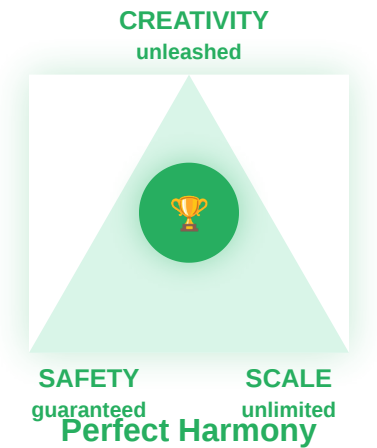
## The Original Challenge

CREATIVITY
vs

⚠️

SAFETY          SCALE
vs               vs

**Pick Any Two**

→

## Our Progress So Far

CREATIVITY
improving

📍

SAFETY          SCALE
strengthening   growing

**Making Progress**

🚀

## The Journey Continues

Our destination: achieving all three in perfect harmony

## Our Ultimate Goal

CREATIVITY
unleashed

🏆

SAFETY          SCALE
guaranteed      unlimited

**Perfect Harmony**

# Three Principles of Lightweight Alignment

A framework we're exploring for building AI that works with humans

## Alignment as Product Design

We're exploring alignment as user experience design with safety constraints. The best solutions might emerge when we design for human needs first.

*"How do users actually interact with this? What behavior signals tell us it's working?"*

## Behavior Over Preferences

Users show us what works through their actions, not their words. Behavior signals—completion rates, engagement patterns, usage flows—may tell us more than surveys.

*"Children vote with their attention. Completion rates matter more than survey responses."*

## Constraints Enable Creativity

We're testing whether the right guardrails can guide expression toward more meaningful outcomes rather than limiting it. Structure might become the foundation for innovation.

*"60-70% constraint level = peak creativity. Structure channels imagination productively."*

# Thank You

Questions? Let's discuss the journey ahead

🚀 📊 🎨

## Vijay Chakilam

**Founder, Hello Kooper**

@thankrandomness          linkedin.com/in/vijaychakilam

This research represents early findings from our AI storytelling platform.
We're excited to continue exploring how constraints can enable creativity at scale.