

Integrated Workflow Demo

Module 7 of 7

Presenter: Vương Cường — Research Assistant, BAI Lab

Complete Research Pipeline Integration

Goal: Demonstrate all tools working together in a real research workflow

What We'll Build

- **End-to-end research project** using base-research-repo
- **All 6 tools integrated** - Git, UV, Pre-commit, HF Hub, W&B, Telegram
- **Professional workflow** suitable for publication-quality research
- **Team collaboration** ready for multi-researcher projects

Why This Integration Matters

- **Systematic research methodology** - No more ad-hoc experimentation
- **Reproducible results** - Complete provenance tracking
- **Professional quality** - Industry-standard development practices
- **Collaborative efficiency** - Team-wide consistency and transparency

Comparison: Before vs After Integration

Aspect	Before Integration	After Integration
Environment Setup	Manual pip installs, environment conflicts	UV: Instant, consistent environments
Code Quality	Inconsistent style, manual reviews	Pre-commit: Automated quality gates
Version Control	Chaotic file management	Git: Professional branching workflow
Data Storage	Email attachments, Drive folders	HF Hub: Versioned dataset management
Experiment Tracking	Spreadsheets, print statements	W&B: Professional experiment logging

Business Impact

- **Faster onboarding** - New researchers productive in minutes
- **Reduced errors** - Automated quality and consistency checks
- **Better collaboration** - Shared tools and standards
- **Publication ready** - Professional research methodology

Step 1: Repository Setup with All Tools

```
# 1. Clone and setup base research repository
git clone https://github.com/bailab/base-research-repo
cd base-research-repo

# 2. Create feature branch for new experiment
git checkout develop
git pull origin develop
git checkout -b feature/sentiment-analysis-experiment

# 3. Setup environment with UV
uv venv
source .venv/bin/activate # macOS/Linux
uv pip install -e .

# 4. Install development dependencies
uv pip install wandb huggingface_hub requests python-dotenv

# 5. Setup pre-commit hooks
uvx pre-commit install

# 6. Setup authentication for external services
huggingface-cli login # Enter HF token
wandb login # Enter W&B API key

# 7. Setup environment variables
cat > .env << 'EOF'
BOT_TOKEN=your_telegram_bot_token
CHAT_ID=your_telegram_chat_id
WANDB_PROJECT=bailab-sentiment-analysis
HF_DATASET_REPO=bailab/sentiment-dataset-v1
HF_MODEL_REPO=bailab/sentiment-model-v1
EOF
```

Step 2: Create Integrated Training Script

```
# src/training/integrated_train.py
import os
import wandb
import requests
from datetime import datetime
from pathlib import Path
from huggingface_hub import upload_folder, download_folder
from datasets import load_dataset
from dotenv import load_dotenv

load_dotenv()

def notify_telegram(message: str):
    """Send Telegram notification."""
    bot_token = os.environ.get("BOT_TOKEN")
    chat_id = os.environ.get("CHAT_ID")

    if bot_token and chat_id:
        url = f"https://api.telegram.org/bot{bot_token}/sendMessage"
        requests.post(url, data={
            "chat_id": chat_id,
            "text": message,
            "parse_mode": "Markdown"
        })

def complete_research_pipeline():
    """Demonstrate complete integrated workflow."""
    try:
        # Initialize W&B experiment
        wandb.init(
            project=os.environ.get("WANDB_PROJECT", "bailab-research"),
            name=f"sentiment-analysis-(datetime.now().strftime('%Y%m%d-%H%M'))",
            config={
                "model_type": "transformer",
                "learning_rate": 2e-5,
                "batch_size": 16,
                "epochs": 3,
                "dataset_version": "v1.0"
            }
        )

        # Telegram: Start notification
        notify_telegram(f"🚀 Starting experiment: {wandb.run.name}")

        # Load data from Hugging Face Hub
        dataset = load_dataset(os.environ.get("HF_DATASET_REPO"))
        notify_telegram(f"📁 Dataset loaded: {len(dataset['train'])} training samples")

        # Simulate training loop
        for epoch in range(wandb.config.epochs):
            # Simulate training metrics
            train_loss = 1.0 - (epoch * 0.2)
            val_accuracy = 0.6 + (epoch * 0.15)

            # Log to W&B
            wandb.log({
                "epoch": epoch,
                "train/loss": train_loss,
                "val/accuracy": val_accuracy,
                "learning_rate": wandb.config.learning_rate
            })

            # Telegram progress update
            if epoch % 1 == 0: # Every epoch for demo
                notify_telegram(
                    f"🔄 **Epoch {epoch + 1}/{wandb.config.epochs}**\n"
                    f"📉 Loss: {train_loss:.3f}\n"
                    f"📈 Accuracy: {val_accuracy:.3f}\n"
                    f"🔗 [View Progress]({wandb.run.url})"
                )

            print(f"Epoch {epoch}: loss={train_loss:.3f}, acc={val_accuracy:.3f}")

        # Save model artifacts
        model_dir = Path("./trained_model")
        model_dir.mkdir(exist_ok=True)

        # Create dummy model files
        (model_dir / "config.json").write_text('{"model_type": "transformer"}')
        (model_dir / "pytorch_model.bin").write_text("dummy_model_weights")

        # Upload to Hugging Face Hub
        upload_folder(
            folder_path=str(model_dir),
            repo_id=os.environ.get("HF_MODEL_REPO"),
            repo_type="model",
            commit_message=f"Upload model from experiment {wandb.run.name}"
        )

        # Final notifications
        notify_telegram(
            f"🏁 **Experiment Complete!**\n"
            f"📊 Final Accuracy: {val_accuracy:.3f}\n"
            f"🔗 W&B: {wandb.run.url}\n"
            f"📄 Model: https://huggingface.co/{os.environ.get('HF_MODEL_REPO')}"
        )

        wandb.finish()
    except Exception as e:
        notify_telegram(f"❌ Experiment Failed: {e}")
        raise
    if __name__ == "__main__":
        complete_research_pipeline()
```

Step 3: Execute Complete Pipeline

```
# Run the integrated training script
python src/training/integrated_train.py

# What happens automatically:
# 1. Pre-commit runs on any code changes
# 2. W&B experiment starts and logs metrics
# 3. Telegram sends real-time updates
# 4. Dataset loads from HF Hub
# 5. Model uploads to HF Hub
# 6. Complete experiment tracking and notifications
```


Step 4: Commit and Create Pull Request

```
# Stage all changes
git add .

# Pre-commit automatically runs:
# - Ruff linting and formatting
# - MyPy type checking
# - File hygiene checks

# Commit with conventional format
git commit -m "feat(training): add complete integrated research pipeline

- Implement end-to-end training with all tools
- Add W&B experiment tracking
- Include Telegram notifications
- Setup HF Hub integration for datasets and models
- Demonstrate professional research workflow

Closes #123"

# Push feature branch
git push -u origin feature/sentiment-analysis-experiment







# Create Pull Request on GitHub
# - Links to issue #123
# - Includes experiment results and W&B dashboard
# - Shows Telegram notification history
# - Demonstrates code quality via pre-commit
```

Step 5: Team Review and Collaboration

```
# Team members can:  
# 1. Review code changes in PR  
# 2. Access W&B experiment dashboard  
# 3. Download trained model from HF Hub  
# 4. Receive Telegram notifications in group chat  
# 5. Reproduce experiment with exact environment  
  
# Reviewer workflow:  
git checkout feature/sentiment-analysis-experiment  
uv venv && source .venv/bin/activate  
uv pip install -e .  
python src/training/integrated_train.py # Reproduces experiment
```

Professional Research Methodology

Reproducibility Checklist

-  **Environment versioning** - UV lock files ensure consistent dependencies
-  **Code versioning** - Git tracks all changes with commit history
-  **Data versioning** - HF Hub provides dataset version control
-  **Experiment versioning** - W&B logs complete experiment provenance
-  **Quality assurance** - Pre-commit prevents low-quality code
-  **Real-time monitoring** - Telegram enables immediate issue detection

Team Collaboration Benefits






- **Consistent development environment** - All team members use same tools
- **Shared experiment visibility** - W&B dashboards accessible to all
- **Standardized code quality** - Pre-commit enforces team-wide standards
- **Centralized data management** - HF Hub provides single source of truth
- **Communication automation** - Telegram keeps everyone informed

Research Quality Improvements

Before Integration: Common Problems

- ✗ "Works on my machine" syndrome
- ✗ Lost experiments and forgotten hyperparameters
- ✗ Inconsistent code quality across team
- ✗ Data scattered across email and drives
- ✗ Manual monitoring leads to missed failures

After Integration: Professional Standards

-  **Systematic experimentation** with complete tracking
-  **Reproducible results** through proper versioning
-  **Publication-ready code** with automated quality control
-  **Efficient collaboration** via shared tools and standards
-  **Proactive monitoring** with real-time notifications

Measurable Impact

- **Setup time:** Hours → Minutes (10x improvement)
- **Experiment reproducibility:** 30% → 95% success rate
- **Code quality issues:** Reduced by 80% via automation
- **Team onboarding:** Days → Hours (5x faster)
- **Research velocity:** More time on science, less on DevOps

Complete Tool Integration

- ✓ **Git** - Professional version control and collaboration
- ✓ **UV** - Fast, reliable Python environment management
- ✓ **Pre-commit** - Automated code quality enforcement
- ✓ **HF Hub** - Proper dataset and model versioning
- ✓ **W&B** - Professional experiment tracking and visualization
- ✓ **Telegram** - Real-time monitoring and team communication

Key Achievements

1. **Eliminated manual DevOps overhead** - Focus on research, not tooling
2. **Established professional standards** - Publication-quality methodology
3. **Enabled true reproducibility** - Complete experiment provenance
4. **Optimized team collaboration** - Shared tools and workflows
5. **Automated quality assurance** - Prevent technical debt accumulation

Research Impact

- **Faster hypothesis testing** - Reduced friction in experimentation
- **Higher quality results** - Systematic methodology prevents errors
- **Better collaboration** - Team-wide visibility and standards
- **Publication readiness** - Professional code and documentation
- **Sustainable practices** - Long-term maintainable research projects

Call to Action

1. **Adopt this workflow** for your current research project
2. **Setup team-wide standards** using these tools
3. **Train new researchers** on this integrated approach
4. **Contribute improvements** back to base-research-repo
5. **Share success stories** with the research community

Transform your research from ad-hoc experimentation to systematic science!