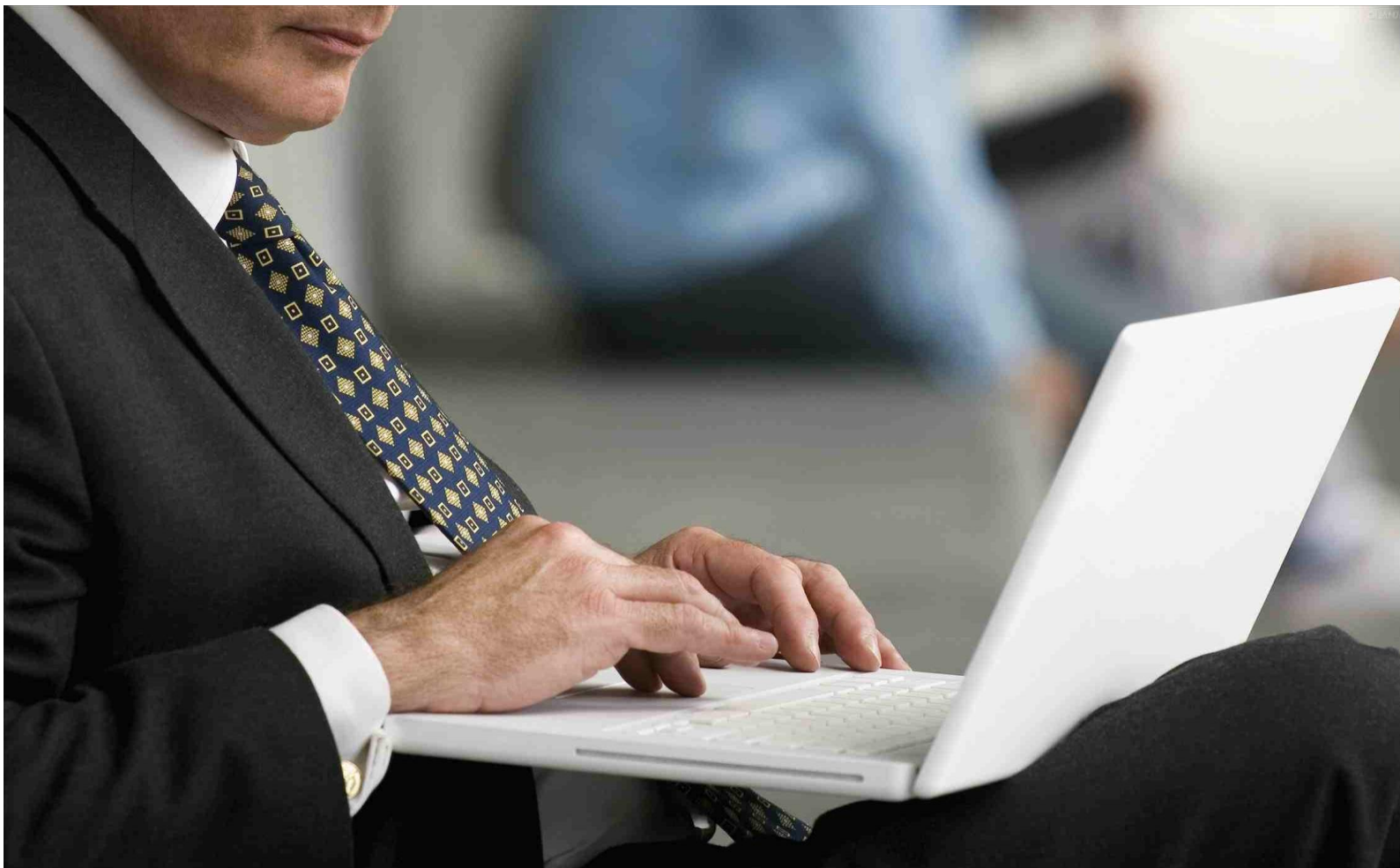


# 数据集成建模指南

在许多组织中，都有大量的数据集成流程重复存在，主要原因之一就是没有可视化的方法可以“看到”当前存在的数据集成过程以及还需要什么过程。

——Anthony David Giordano

- 软件开发生命周期内的数据集成建模
- 利用流程建模进行数据集成
- 解读数据集成建模中的数据模型
- 数据集成模型开发工具指南



# 数据集成建模指南

*数据集成流程的开发类似于数据库的开发。在开发数据库过程中，业务需求的蓝图或者模型必须确保对需要的部分有清楚的理解。——Anthony David Giordano*

# 结构化数据集成建模与数据集成架构



本此数据库电子书节选自《数据集成蓝图和建模》一书，读者可以了解到如何为新的数据集成设计流程构建业务案例，以及如何为数据集成建模改进开发流程。读者还可以获得为数据集成和设计数据集成架构模型利用流程建模的技巧，还会了解到三种数据

集成建模类型——物理建模，逻辑建模和概念建模。

第一部分主要介绍了一种新的设计技术，它是用来分析和设计数据集成流程的。这种技术使用图形化流程建模的数据集成视图，类似于为数据模型提供的实体关系图那样的图形化视图。

## 新设计流程的业务案例

对于数据集成流程的大规模复制问题，有一个情形如下：

如果你没有看到某个过程，你就会重复该过程。

在许多组织中，都有大量的数据集成流程重复存在，主要原因之一就是没有可视化的方法可以“看到”当前存在的数据集成过程以及还需要什么过程。这与曾经给数据建模规程带来麻烦的问题很类似。

在 20 世纪 80 年代早期，许多组织都在大规模复制客户和交易数据。这些组

织看不到他们数据环境和大规模复制的“全貌”。一旦组织开始文档化记录并充分利用实体关系图(数据模型的可视化展现形式)，他们就能看到大量复制工作，而且增加不必要的复制会降低现存表的复用度。

数据集成流程的开发类似于数据库的开发。在开发数据库过程中，业务需求的蓝图或者模型必须确保对需要的部分有清楚的理解。在数据集成的案例中，数据集成设计者和数据集成开发人员需要该蓝图或项目工件，来确保关于需要移动数据的源、转换和目标的业务需求已经通过共同一致的方法进行了清晰的交流。对专门为数据集成设计流程模型的使用将实现该需求。

图 1 描述了项目中需要的数据模型类型，你可以看到它们与为数据集成开发的模型有多相似。



图 1 建模示例：数据和数据集成。



在大部分项目中，分析，设计以及构建 ETL 或者数据集成过程的通用方法都涉及到数据分析文档化需求，在微软 Excel 数据表中定义源到目标的映射。这些数据表被提供给 ETL 开发者，用来设计和开发映射，图表以及开发源代码。

手工把源系统和目标系统的集成需求文档化记录到像 Excel 这样的工具中，然后把它们再映射到 ETL 或者数据集成包，事实证明这种做法非常耗时，而且容易出错。例如：

**消耗的时间。**从源系统向 Excel 数据表复制数据会花费相当多的时间。相同源的信息必须在 ETL 工具中重新生成键值。此源和在 Excel 中收集的目标元数据基本很难复用，除非有大量的人工审查和维护流程。

**非值增加分析。**利用转换需求捕获源到目标映射包含有价值的导航性元数据，这些数据可以被用于数据发展分析。在 Excel 电子表格中捕获这种信息，不能提供清晰自动的方法捕获这种有价值的信息。

**映射错误。**尽管我们付出最大的努力，但是手工操作数据通常还是会出错的，例如，在 Excel 电子表中可能会将“INT”数据类型转成“VARCHAR”数据类型，这需要数据集成设计者花时间分析和纠正。

**缺少标准：不一致的详细程度。**执行源到目标映射的数据分析师们很容易以不同的完整程度来捕获源/转换/目标需求，这取决于该分析师的技能和经验。一旦在需求和数据集成流程设计的详细程度上出现了不一致，就可能会给开发人员

阅读源到目标映射的文档(通常是 Excel)时造成误解，这样就会导致编码错误和浪费时间。

**缺乏标准：文件格式不一致。**大部分环境对不同的文件格式有多种方式提取。工作的重点和方向必须是朝着一次读取，多次写入的概念进行，同时要保持抽取、数据质量、转换以及加载格式的一致性。

要提升数据集成流程的设计和开发效率，还有时间、一致性、质量以及可重用性，对于开发数据模型采用同样严密的数据集成使用一套图形化建模设计技术是很有必要的。

对于诸如数据集成这类信息技术应用,流程建模是一种经过尝试并证明可行的方法。通过对数据集成应用流程建模技术，虚拟化和标准化的问题也会涉及到。首先，我们来了解一下流程建模的类型吧。

## **利用流程建模进行数据集成**

流程建模是以某种详细程度展示系统相关流程的一种手段，利用指定类型的图表通过一系列流程展示数据流。流程建模技术通常图形化地展示具体流程，以便更清晰地理解，交流并在设计和开发系统流程中的涉众之间进一步精细化。

流程建模不像数据建模，有几种不同的流程模型类型，它与不同的流程交互类型有关。这些不同的模型类型包括流程依赖图，结构层次图以及数据流图。数

据流图是这些流程模型类型中最广为人知的一种，它可以进一步被提炼成几种不同的数据流图类型，比如背景图，还有代表不同层次和流程类型数据量的 0 级图、1 级图以及“叶子级”图。

数据集成建模是一种流程建模技术，关注把工程化数据集成流程转化成普通数据集成架构。

通过利用不同层次的概念以及流程建模的类型，我们已经为数据集成流程开发出了一套流程建模方法，描述如下：

## 数据集成建模概览

数据集成建模是一种技术，考虑到了基于数据集成架构需求类型的模型类型和基于系统开发生命周期(SDLC)的模型类型。

## 为数据集成架构建模

流程建模类型或者数据集成模型类型依赖于数据集成参考架构中需要的流程化类型。通过把参考架构用作框架，我们能为离散数据集成流程和着陆区创建具体流程模型类型，如图 2 所示。

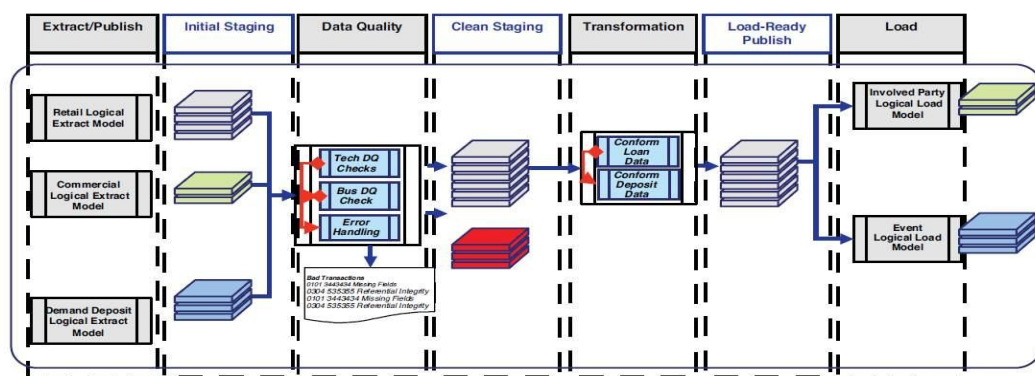




图 2 为架构设计模型。

总之，这些离散数据集成层变成了流程模型类型，形成了完整的数据集成流程。目标是开发一种技术，可以引导设计者基于一组通用的流程类型为数据集成流程建模。

### 软件开发生命周期内的数据集成建模

数据集成模型遵从与软件开发生命周期中数据建模时出现的需求和设计抽象精炼通用的级别。正如存在概念的，逻辑的和物理的数据模型，也存在概念的，逻辑的和物理的数据集成需求，需要在软件开发生命周期的不同点进行捕获，它们可能在流程模型中有所展现。

下面是每种模型类型的简要说明，关于角色、步骤以及模型示例的更完整定义将会在本章的后面进行阐述。

**概念数据集成模型定义。**为目标系统产生一种无需实施的数据集成需求展现，将作为确定他们怎样能得到满足的基础。

**逻辑数据集成模型定义。**在数据集层面产生详细的数据集成需求展现，详细描述了转换规则和目标逻辑数据集。这些模型仍然被认为是技术无关的。在逻辑层面的重点是在于真正源表以及建议目标存储的捕获。

**物理数据集成模型定义。**在组件层级产生数据集成规格的详细描述。它们应

该被以基于组件的方法展现，而且能展现数据将如何在选定的开发技术中通过数据集成环境优化流程。

### 在参考架构之上构建模型

在系统开发生命周期中构建数据模型是相对容易的过程。一个概念数据模型通常只有一个逻辑模型，一个逻辑数据模型通常也只有一个物理模型。尽管一个模型内的实体可以被进一步分解或者规范化，事实上也很少会需要把一个数据模型分成两个独立的模型。

过程模型历来会被进一步划分成离散的功能。例如，在图 3 中，数据流图的顶部流程是一个内外关系图，它被进一步分成几个独立的功能模型。

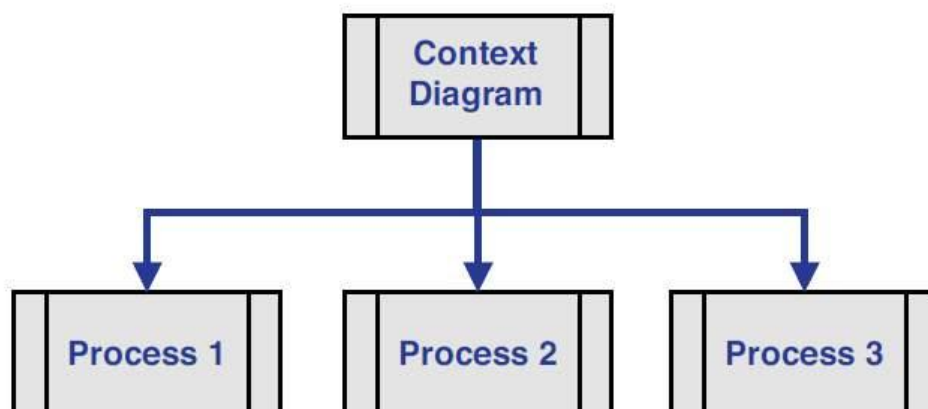


图 3 传统的流程模型：数据流图。

数据集成模型也被分成功能模型，是基于数据集成参考架构和系统开发生命周期状态进行的。

图 4 描绘了概念的，逻辑的和物理的数据集成模型是如何分解的。

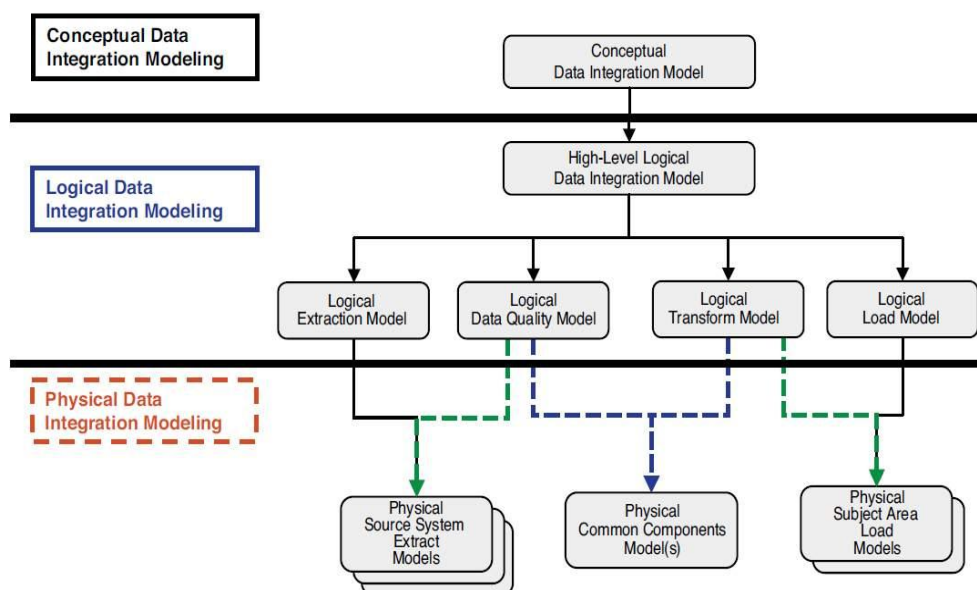


图 4 系统开发生命周期的数据集成模型。

(作者: Anthony David Giordano 译者: 冯昀晖 来源: TT 中国)

# 解读数据集成建模中的数据模型

在本部分中，读者将了解到概念数据集成模型和逻辑数据集成模型，包括高级逻辑集成模型。读者还将了解到把逻辑数据模型转换成物理数据集成模型的相关技巧。此外，我们还能了解到关于利用基于目标的数据集成设计技术和关于设计提取核查过程信息的概要介绍。

## 概念数据集成模型

概念数据集成模型是针对目标系统数据集成需求的一种无需实现的展示，该系统将作为基本的“范围”，定义了它们要如何才能被满足，同时也是为了源系统分析，任务和持续时间以及资源的项目规划的目的。

在这个阶段，只需要确定主要的概念性流程，充分理解用户数据集成并规划下一阶段需求。

图 5 提供了概念性数据集成模型的一个示例。

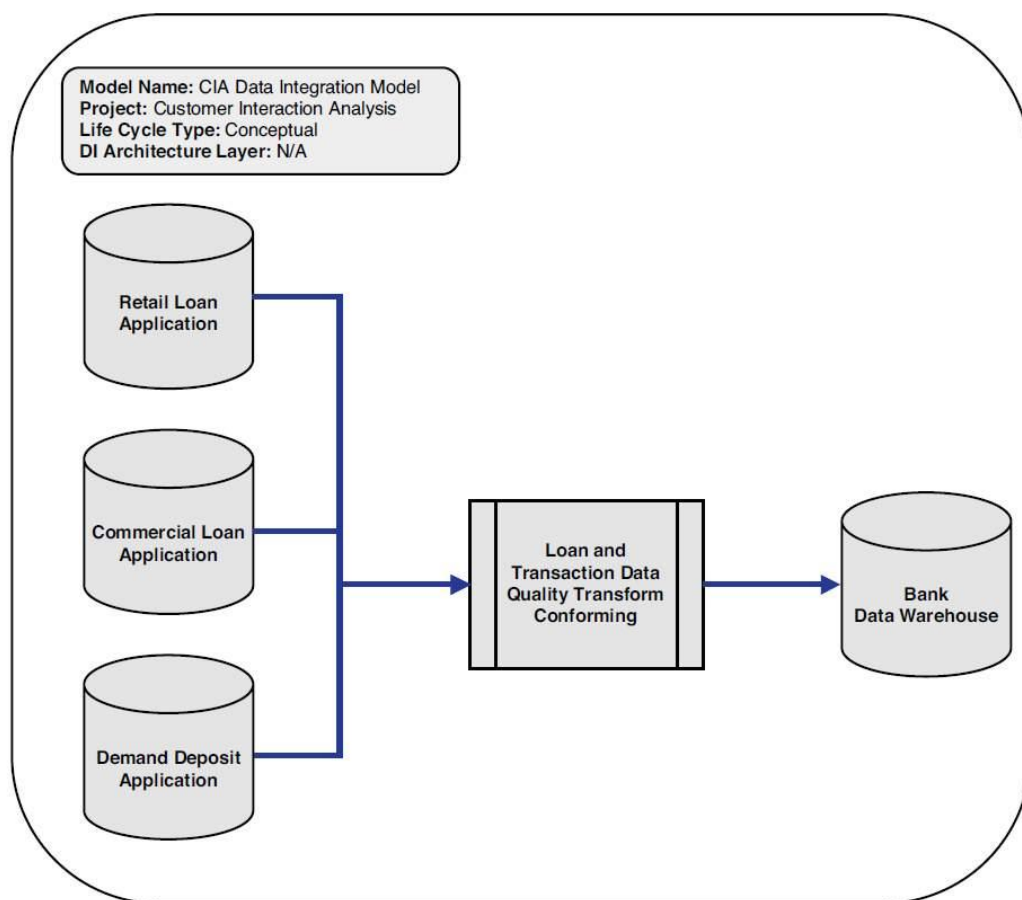


图 5 概念性数据集成模型示例。

## 逻辑数据集成模型

逻辑数据集成模型对捕获首次过滤源映射、业务规则、目标数据集(表或者文件)的数据集成需求产生一组详细描述。这些模型为预期的数据集成应用描绘逻辑提取，数据质量，转换，以及加载需求。这些模型仍然被认为是与技术无关的。本文后面的部分讨论了各种逻辑数据集成模型。

## 高层次逻辑数据集成模型

高层次逻辑数据集成模型定义了项目和系统的范围和界限，该模型通常是从



概念数据集成模型中延伸而来的。高层次数据集成图表提供的规则与为数据流图提供的上下文图表一样。

如下图 6 所示，高层次逻辑数据集成模型提供了数据集成系统需要的结构，还提供了逻辑模型的提纲，比如抽取、数据质量、转换以及加载组件。

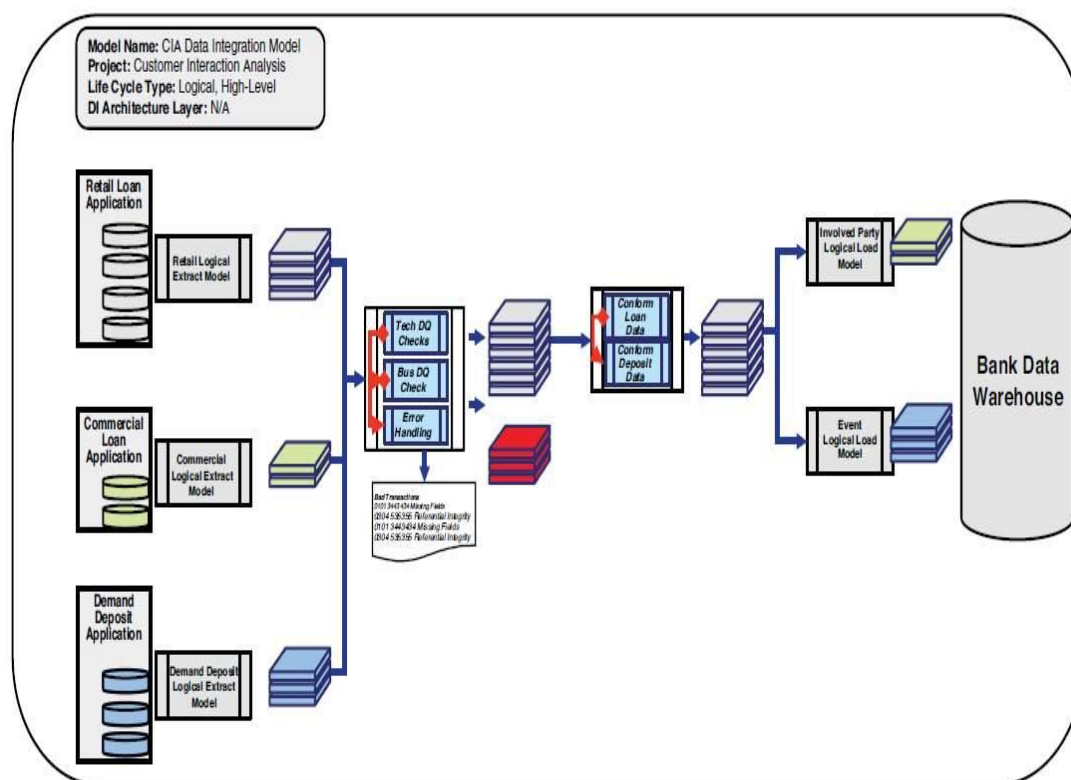


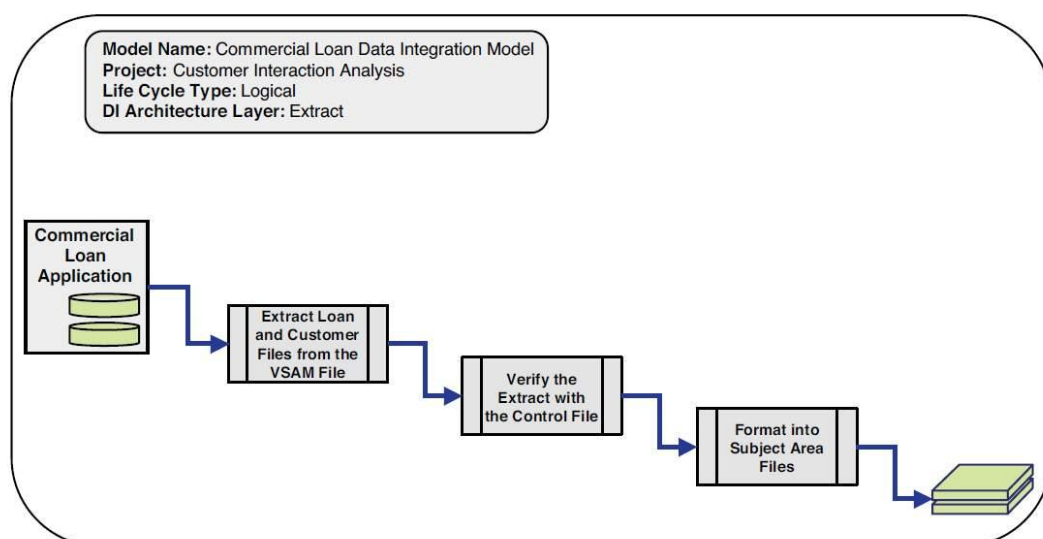
图 6 逻辑高层次数据集成模型示例。

## 逻辑抽取数据集成模型

逻辑抽取数据集成模型决定哪些主题领域将需要从源抽取出来，比如：哪些应用、数据库，平面文件以及非结构化的数据源。

源文件格式应该被映射成属性/字段列/域层次。一旦提取了，源数据文件应该被默认为初始临时区域加载。

图 7 显示了一个逻辑提取模型。



**图 7 逻辑提取数据集成模型示例。**

提取数据集成模型包含两个独立的子过程或者组件：

从源系统中取出数据。无论数据实际上是从源系统中提取的，还是从消息队列或平面文件中捕获的，指向源的网络连接必须是确定的，表或文件的数量必须被审查，而且要提取的文件以及以什么顺序提取它们必须是确定的。

把数据格式化为主题领域文件。主题领域文件提供了从源到最终目标区域的封装层。提取数据集成模型的第二个主要组件是从源格式到通用主题域文件格式的梳理，例如：把西贝尔客户关系管理软件的一组表映射到客户的主题领域文件

中。

## 逻辑数据质量数据集成模型

逻辑数据质量数据集成模型针对预期数据集成流程包含业务和技术数据质量检查点，如图 8 所示。

不管技术还是业务的数据质量需求，每种数据质量数据集成模型都应该包含生成清洗文件，拒绝文件和拒绝报告(在选定的数据集成技术中会举例说明)的能力。

同样，整个数据集成流程的错误处理也应该设计为可重用的组件。

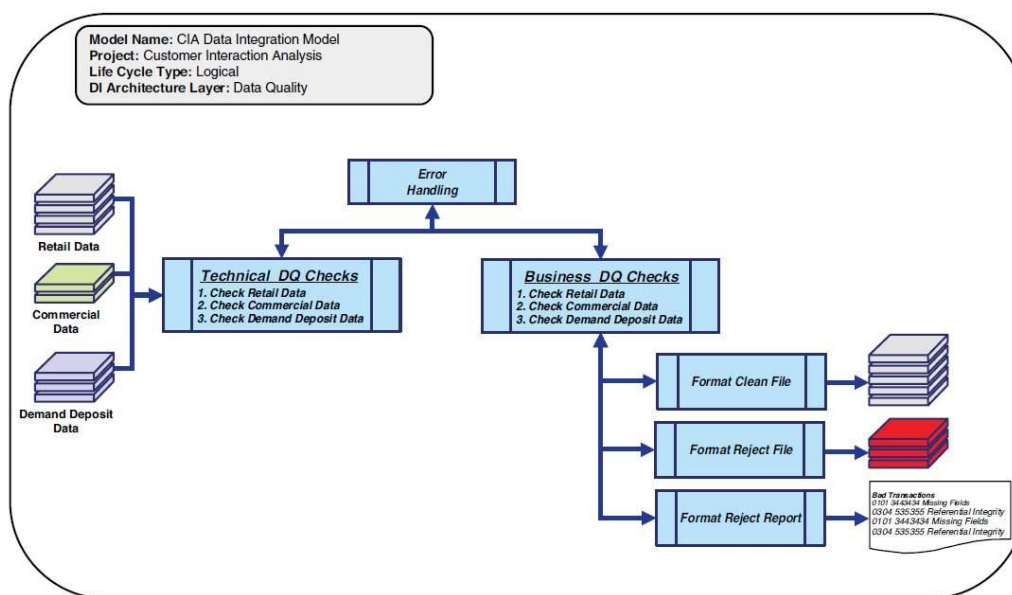


图 8 逻辑数据质量数据集成模型示例。

正如在第二章的数据质量架构流程中讨论的，明确的数据质量流程会产生清

洗文件，拒绝文件和拒绝报告。基于组织的数据治理过程，拒绝文件可以在手工或自动的流程再造利用到。

## 逻辑转换数据集成模型

逻辑转换数据集成模型认定了一个逻辑水平，转换(关于计算、分割、处理和浓缩)需要被在提取的数据上执行以满足关于聚合、计算和构造的商业智能需求。

如图 9 所示。

在转换流程中定义的转换类型是由适应、计算和聚合数据成企业信息业务需求决定的。

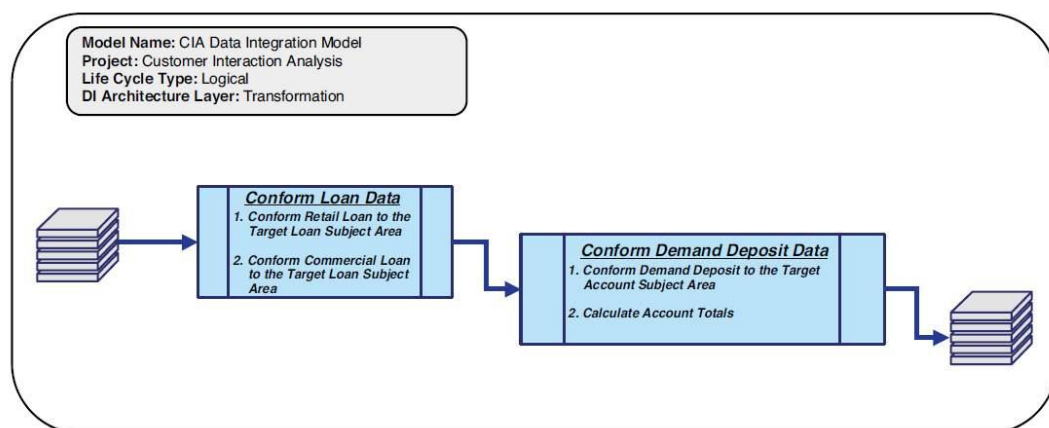


图 9 逻辑转换数据集成模型示例。

## 逻辑加载数据集成模型

逻辑加载数据集成模型确定在一个合理的水平，需要加载转换和清洗后的数据到主题领域的目标数据资产库中。如图 10 所示。

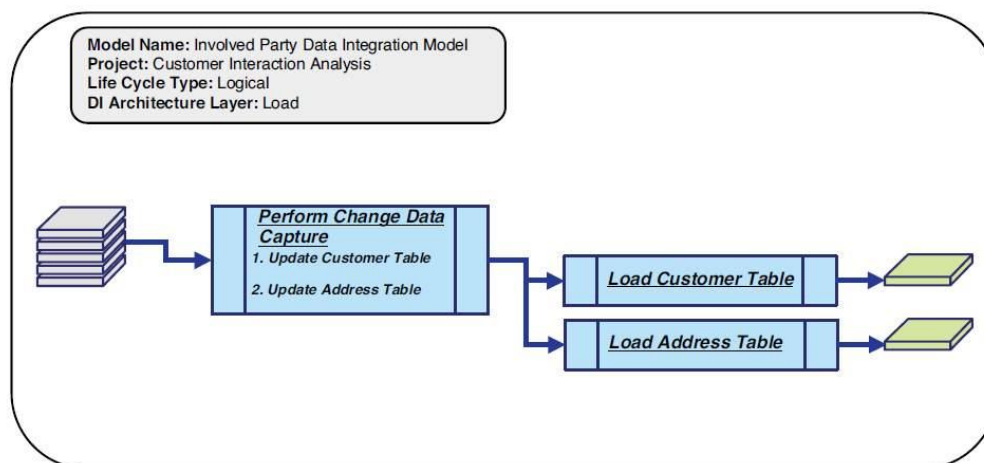


图 10 逻辑加载数据集成模型示例。

在既定的目标数据库中根据目标和主题领域涉及加载流程允许对子过程进行定义，进一步从源数据在目标中封装变化，防止重大维护。一个典型的例子就是，如果出现了数据库物理结构变化，只有主题领域加载任务需要修改，对提取和转换流程影响很微小。

### 物理数据集成模型

物理数据集成模型的目的是在目标数据集成技术中在组件层产生一个详细的数据集成规范表述。

在物理数据集成建模中有一个主要的思想就是要判断怎样最好地形成逻辑设计并应用能优化性能的设计技术。

### 转换逻辑数据集成模型为物理数据集成模型

正如在数据建模中有一个从逻辑到物理数据模型的转换，在数据集成模型中



也存在相同的转换。逻辑数据集成建模判断要提取什么，数据质量，转换以及加载。物理数据集成利用基于目标的设计技术，为如何在物理数据集成模型中设计“怎么做” 确保各种组件能在数据集成环境中最优化执行提供指导。

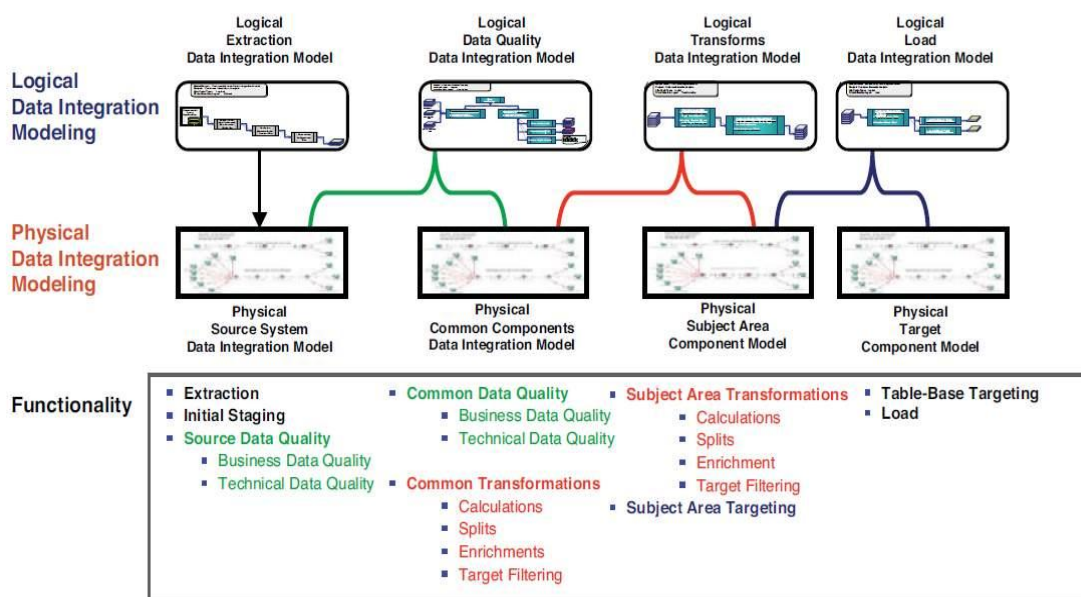
(作者: Anthony David Giordano 译者: 冯昀晖 来源: TT 中国)

# 数据集成建模开发工具

基于目标的数据集成设计技术是基于主题领域加载和位于这些主题领域的源系统创建物理数据集成组件的一种方法。它基于在每种数据集成模型类型中本地与企业级使用数据迁移的模式把逻辑功能分组成可重用的组件。

例如，在大部分数据集成流程中，有源系统级的也有企业级的数据质量检查。基于目标的技术实现的功能，既与将要使用的流程(在这种情况下是提取流程)接近，也在通用组件模型中对企业能力进行了分组。

例如，对于具体源系统的数据质量检查，基于目标的技术简单地把逻辑转移到提取流程，而本地转换被转移到加载流程中，分组企业级数据质量和转换被按照通用组件级别分组。如图 11 显示。



## 图 11 在“是什么”和“如何做”之间分配逻辑功能

基于目标的数据集成设计技术不是一个新的概念：耦合和内聚性，模块化，对象和组件都是把“原材料”分组成可以理解的和高可用工作单元的技术。基于目标的技术是在数据集成模型中模块化核心功能的一种简单方法。

### 物理源系统数据集成模型

提取数据集成模型的源系统从源系统提取数据，执行源系统数据质量检查，然后使数据与具体主题领域文件格式相适应。如图 3.12 所示。

逻辑提取模型与物理源系统数据集成模型的主要差异在于，它关注于最终设计考虑因素，需要从具体源系统提取数据。

### 设计提取验证流程

从源系统文件来的数据被通过控制文件提取并验证。控制文件是一种数据质量检查，验证数据行数和总量控制(举例来说，贷款数量为了针对具体源提取验证被合计起来)。

具体源系统应用的数据质量规则正是在这里。应用具体源系统数据质量规则的基本原理在特定的源系统，而不是在一个整体数据质量任务，这样对维护和性能有好处。巨大的数据质量工作变成了维护的噩梦。它还需要不必要的系统内存量来加载所有数据质量流程和变量，这会减慢整个工作流程的时间。

跨系统依赖在这种模型中应该被处理。例如，连接协议到一起的关联关系在这里应该被处理。

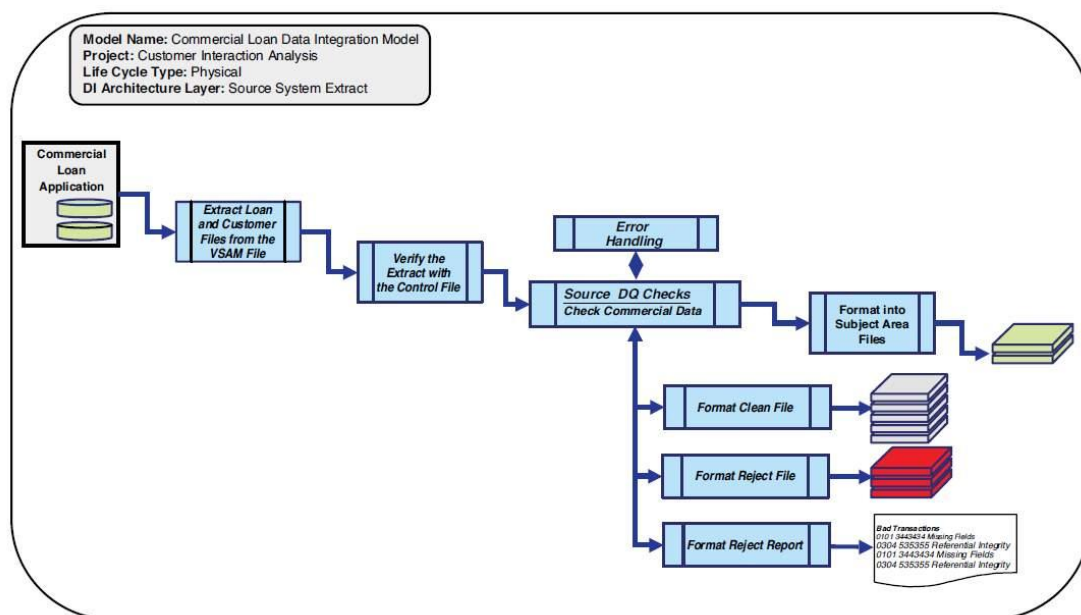


图 12 物理源系统提取数据集成模型示例。

## 物理通用组件数据集成模型

物理通用组件数据集成模型包含企业级业务数据质量规则和通用转换，这些转换将被多种数据集成数据应用。该架构层是整个数据集成流程中重用性的至关重要焦点，尤其强调利用现存的转换组件。任何新组件都必须满足可重用性的标准。

最终，在设计通用组件数据集成模型中，处理流程是在并行被构建的地方检查的，为设计基于预期数据卷并在当前数据集成技术的限制内进行。

## 通用组件数据质量数据集成模型

通用组件数据质量集成模型通常是非常“瘦小”的流程模型(功能较少)，它使用企业级数据质量规则。一般来讲，具体的源系统数据质量规则本质上是技术性的，而业务数据质量规则往往是在企业级应用的。

例如，性别或者邮政编码被视为是可以针对所有待处理数据应用数据质量规则的业务规则。图 13 描绘了通用数据质量数据集成模型的一个示例。

请注意，具体源数据质量规则已经被转移到了物理源系统提取数据集成模型，更精简的数据质量流程是在通用组件级别。更少的数据确保数据流不会收到不必要的限制，而且整体处理性能会得到改善。

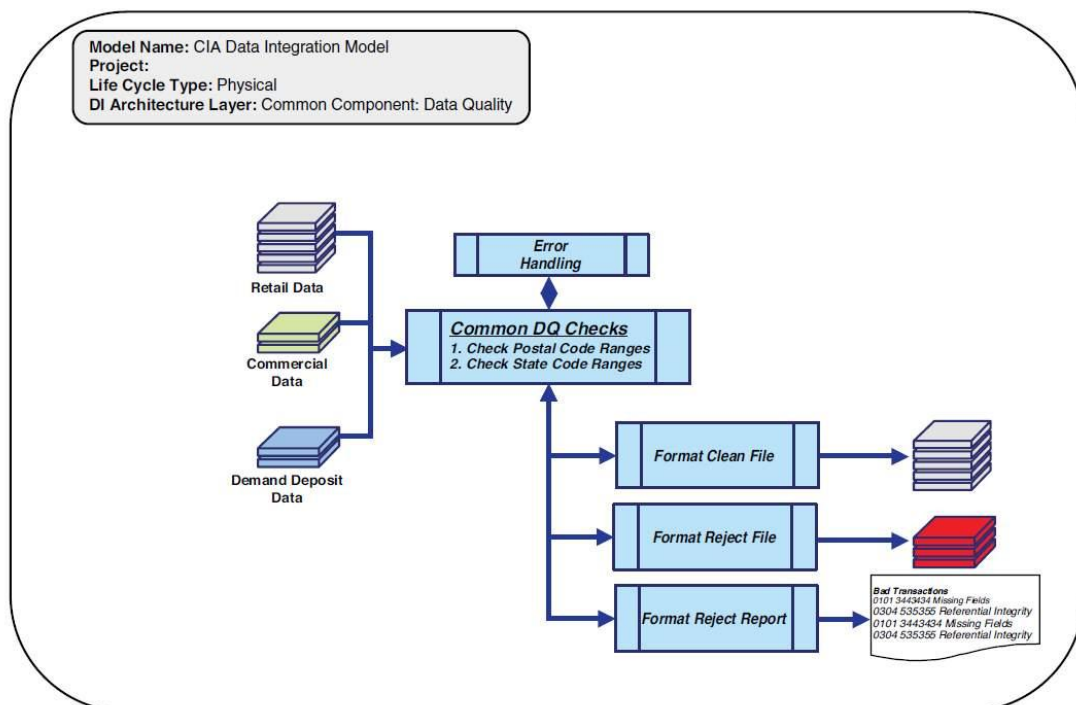


图 13 通用组件——数据质量数据集成模型示例。



## 通用组件转换数据集成模型

最常见的转换是那些使数据符合企业数据模型的转换。需要具体聚合和计算的转换被转移到主题领域加载，或者是转移到它们该去的地方，数据在主题领域被转换。

在企业级聚合和计算方面，通常非常少；大部分转换是针对具体主题领域的。

图 14 描绘了一个通用组件转换数据集成主题领域模型示例。

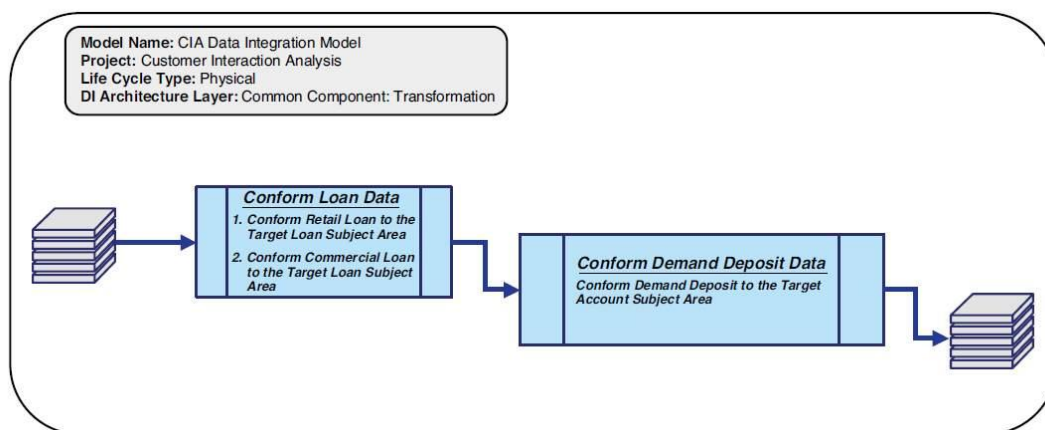


图 14 通用组件——转换数据集成模型示例。

请注意，需求沉淀层的聚合已经从通用组件模型中移除了，已经本着“把功能转移到需要的地方”这一思想转移到主题领域加载了。

## 物理主题领域加载数据集成模型

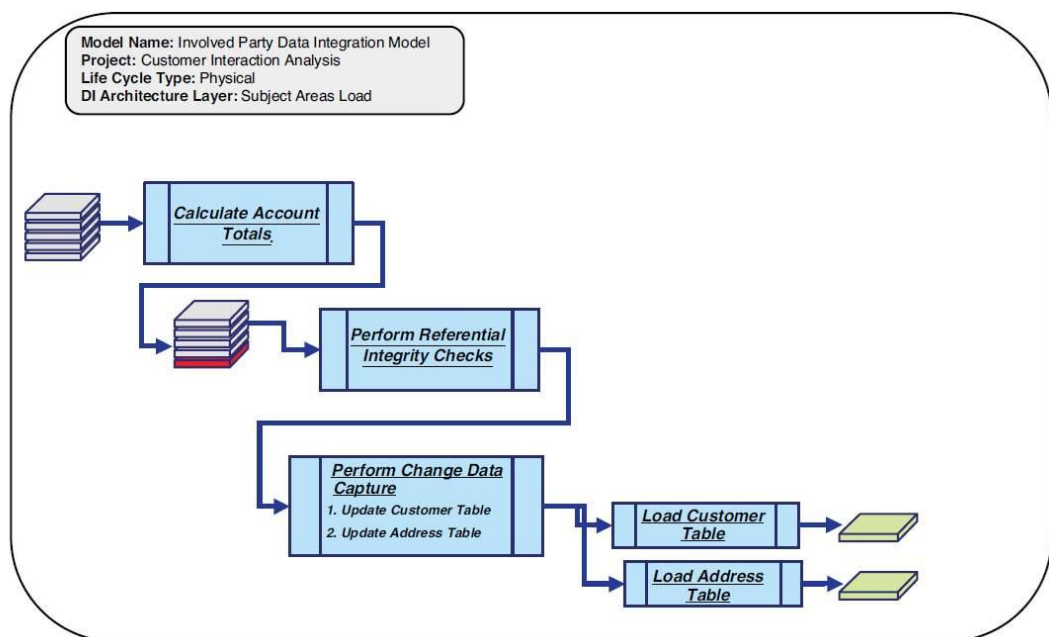
主题领域加载数据集成模型基于主题领域(目标分组)依赖在逻辑上把“目标表”分到一组，并为原系统处理(间接层)简化服务。

加载数据集成模型主题领域执行以下功能：

- 加载数据
- 刷新快照加载
- 执行变更数据捕获

正是在主题领域加载数据集成模型中，主键和外键将被创建，参考完整性被确认，变更数据捕获被处理。

除了为可理解性和可维护性按主题领域分组数据简化，按主题领域分组数据在逻辑上限制了每次处理承载的数据量，因为尽可能少量地承载数据是很重要的，尽管这些流程已经对性能影响做了最小化处理。图 15 展示了物理数据集成主题领域模型的一个示例。



**图 15 物理主题数据领域加载数据集成模型示例。**

### **逻辑数据集成模型 VS. 物理数据集成模型**

在这些工作中总会遇到这样一个问题：“已经需要有一组逻辑数据集成模型了，还需要再有一组物理数据集成模型吗？”

对数据集成模型的答案与数据模型的答案是一样的：“要看情况”。它取决于在他们的元数据管理方面将创建、管理并掌握模型的数据管理组织成熟度，也取决于其他数据管理工件(比如：逻辑数据模型和物理数据模型)。

### **开发数据集成模型的工具**

关于数据集成建模的首要问题之一就是：“你将用什么工具构建模型？” 尽管可以采用像微软公司的 Visio，甚至是微软公司的 PowerPoint 这类图表工具，但是我们仍然主张使用某一款商业数据集成包来设计和构建数据集成模型。

像 Visio 这类绘制图表的工具需要手工创建和维护，才能确保它们与源代码和 Excel 电子表格保持同步。维护方面的代价常常会超过手工创建模型的益处。通过利用数据集成包，现存数据集成设计(比如：提取数据集成模型)可以被审查是否在其它数据模型中有重用的可能性，何时利用，对实际数据集成工作的维护也在模型更新时被执行。另外，通过使用像 Ab Initio，IBM Data Stage 或者 Informatica 这类数据集成包来创建数据集成模型，组织将能进一步利用到它所

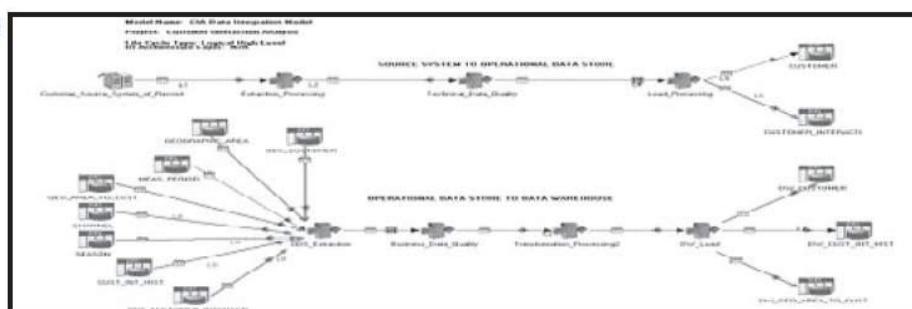
做的技术投资。

图 16 提供了用 Ab Initio , IBM Data Stage 和 Informatica 构建的高层次逻辑数据集成模型示例。

Ab Initio



IBM Data Stage



Informatica



图 16 几个数据集成模型。

在使用数据集成包进行数据集成建模方面的经验表明，数据集成项目和经验中心已经看到了增加提取，转换和本地代码标准化、以及质量的益处。利用数据集成包的关键益处有以下几点：

端对端通信。利用数据集成包有利于更快速地从数据集成设计人员到数据集成开发人员之间传递需求，它们利用相同通用的数据集成元数据。在相同的包中利用相同的元数据从逻辑设计到物理设计转移，加速了传递过程，减少了传输问题和错误。例如，源到目标的数据定义和映射规则不是必须在技术之间传递，因而降低了映射错误。在从逻辑数据模型向物理数据模型转换的数据建模工具中也能发现类似益处。

重点利用企业模型开发。把数据集成需求捕捉为逻辑数据集成模型和物理数据集成模型给组织提供了把这些数据集成模型整合到企业数据集成模型的机会，进一步使信息管理环境变得成熟，增加了整体可复用度。它还提供了重用源提取，目标数据加载和在数据集成软件包元数据引擎中的通用转换的能力。这些物理数据集成任务被保存在相同的元数据引擎中，彼此可以相互链接。他们还可以被链接到其它现存的元数据对象，比如逻辑数据模型和业务功能。

在流程中更早地捕获导航元数据。通过在数据集成软件包中存储逻辑数据集成模型和物理数据集成模型，组织拥有了对单个源或目标任务执行更彻底影响分析的能力。在流程中更早地利用转换需求对源到目标映射元数据的捕获，还增加了在单元测试和系统测试中捕捉映射错误的可能性。此外，由于元数据是被自动捕获的，所以也就更容易捕获和管理。

## **基于行业的数据集成模型**



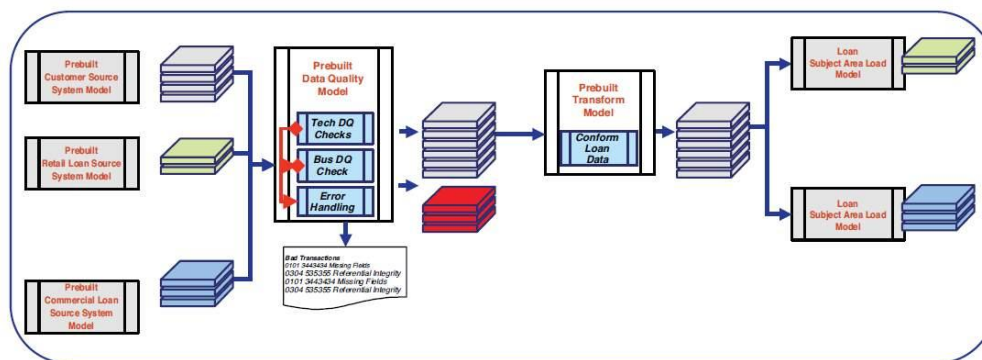
为了降低风险，加快数据仓库项目的设计工作，为数据仓库预构建的数据模型已经由 IBM 公司，甲骨文公司，微软公司和 Teradata 公司开发出来了。

随着数据集成建模概念的成熟，预构建数据集成模型正在那些业界数据仓库数据模型的支持下开发。

预构建数据集成模型把行业数据仓库模型用作目标和知名商业源系统进行提取。有了基于行业的源系统和目标，很容易利用预构建的源到目标映射开发数据集成模型。例如，在银行系统中，就有通用的源系统，比如下面这几种：

- 商业和零售贷款制度。
- 活期存款系统。
- 企业资源系统(比如 SAP 和 Oracle)。

这些已知的应用可以被重新映射到基于行业的数据仓库数据模型。基于实际的项目经验，对基于行业数据集成模型的利用可以极大地削减数据集成项目的时间和成本。图 17 展示了基于行业数据集成模型的一个示例。



**图 17 基于行业的数据集成模型示例。**

在前面的例子中，行业数据集成模型提供如下内容：

- 预构建提取来自客户、零售贷款和商业贷款系统的流程。
- 在目标数据模型中基于已知数据质量需求预构建数据质量流程。
- 基于目标数据模型主题领域预构建加载流程。
- 基于已知数据集成架构、源系统以及目标数据模型从现存设计开始，提供加速数据集成应用开发的框架。

## 总结

数据建模是对数据的一种图形化设计技术。在数据集成工作中，数据集成建模是针对数据集成参考架构利用图形化流程建模技术设计数据集成流程的一种技术。

本此数据库电子书详细介绍了数据集成模型的类型，包括概念建模，逻辑建模和物理建模，介绍了基于数据集成参考架构流程层细分模型的方法以及每种不同逻辑和物理数据集成模型类型的示例。

它涵盖了从逻辑数据集成模型到物理数据集成模型的转换，最恰当的表述可能是如何从“做什么”转向了“怎么做”。

## 我们的编辑团队

您若有何意见与建议，欢迎[与我们的编辑联系](#)。

诚挚感谢以下人员热情参与 TechTarget 中国《数据库电子书》的内容编辑工作！

诚邀更多的数据库专业人士加入我们的内容建设团队！



**沈宏**

TechTarget中国特邀技术编辑。具有丰富的软件开发及测试经验，多年以来一直致力于数据库优化、性能测试等方向的研究与探索。



**冯昀晖**

TechTarget中国特邀技术编辑。资深软件工程师，有超过7年的政府和企业信息化软件解决方案经验，熟悉SQL Server、Oracle等数据库技术，爱好阅读、健身和中国象棋。



**孙瑞**

TechTarget 中国数据库网站编辑，四年网络媒体从业经验。负责“[TT 数据库](#)”网站的内容建设，熟悉数据库以及商业智能等企业信息化领域，拥有计算机学士学位。



关注 [TT 数据库新浪微博](#) 及时了解数据库技术讯息