



整车先进设计制造技术全国重点实验室
State Key Laboratory of Advanced Design and Manufacturing Technology for Vehicle



湖南大学
HUNAN UNIVERSITY

人形机器人遥操作 Humanoid Teleoperation

报告人：刘天适

1.H2O:Human-to-Humanoid CMU



以往工作

基于模型的
控制器

计算成本高

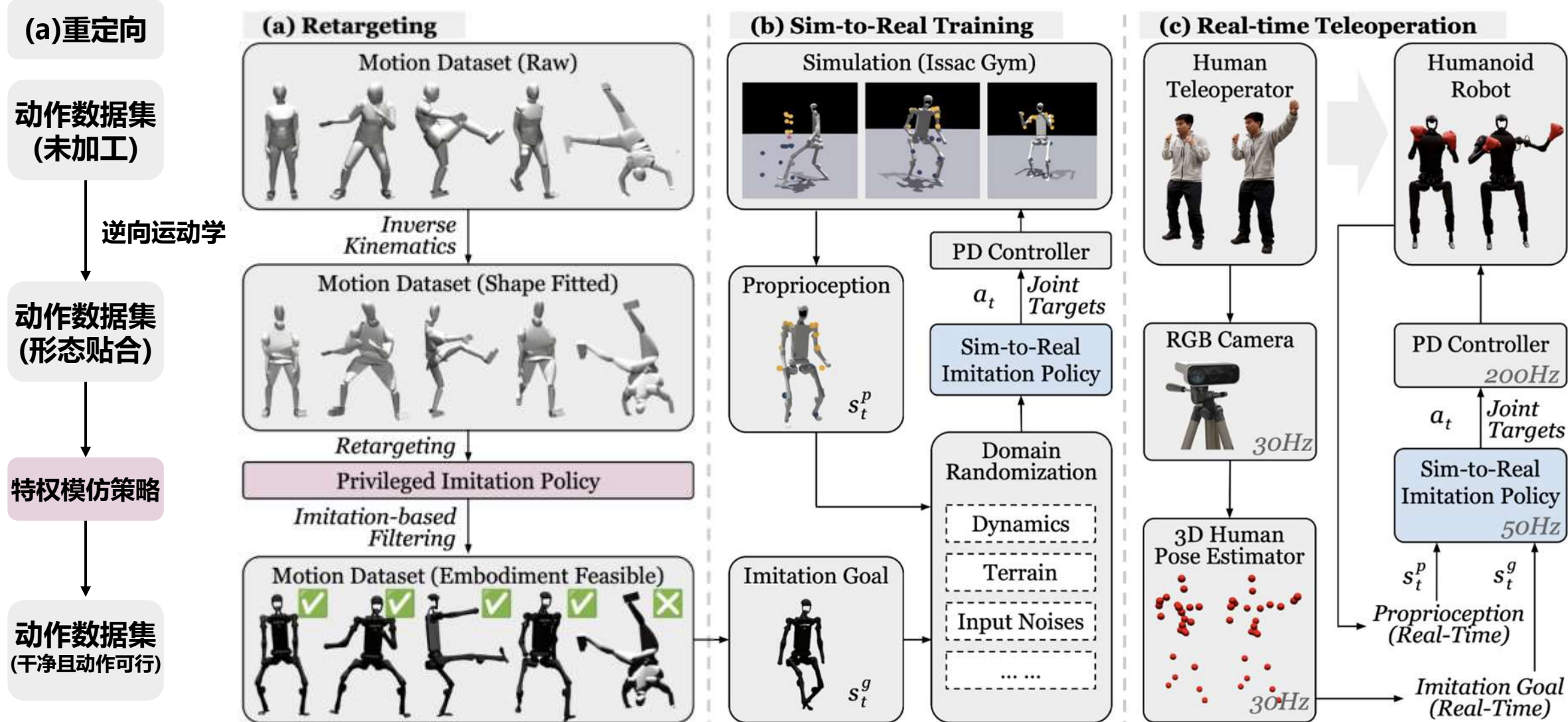
零样本迁移学习

动作捕捉

依赖外部设备
实时性
全身动作的精确模仿

仅用RGB摄像头
实现动作捕捉

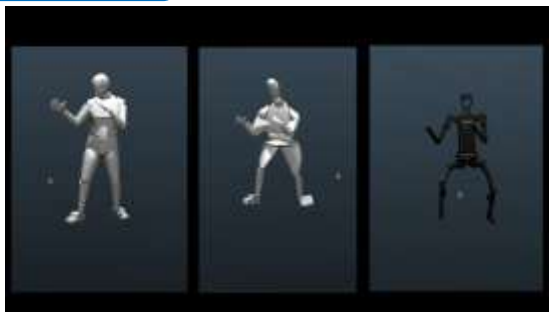
整体框架(包括abc三阶段)



整体框架(包括abc三阶段)

(a)重定向

动作数据集
(未加工)



逆向运动学

调整SMPL身体模型参数
以与机器人结构对齐
(关节对齐, 体型优化, 动作调整)

动作数据集
(形态贴合)

是一种**强化学习**算法, 可以访问完整
刚体状态信息, 包括机器人的全局3D刚
体位置、方向、线速度和角速度等。

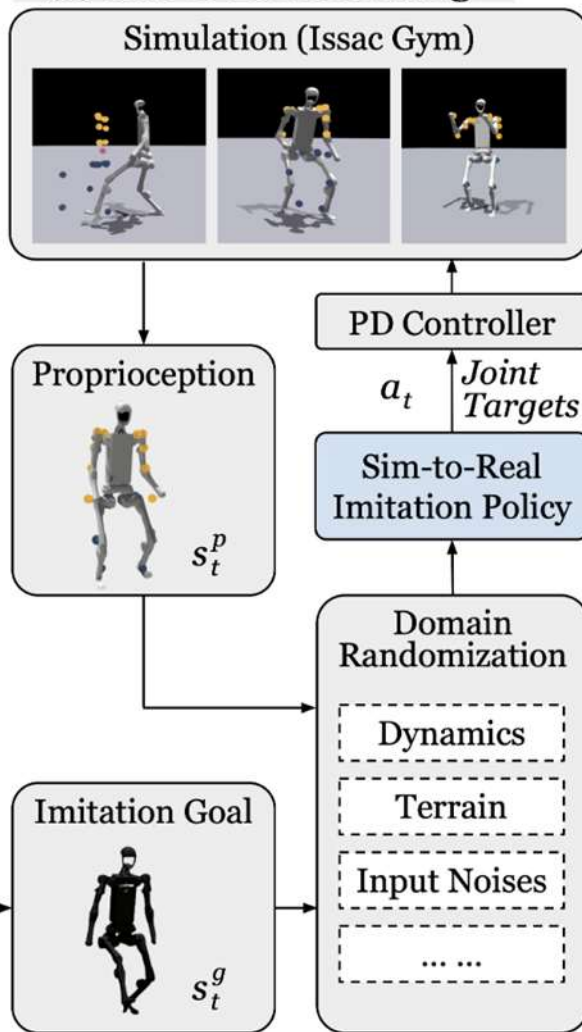
特权模仿策略

这种策略在仿真环境中对机器人进行
训练, 保留机器人**能够执行**的动作。

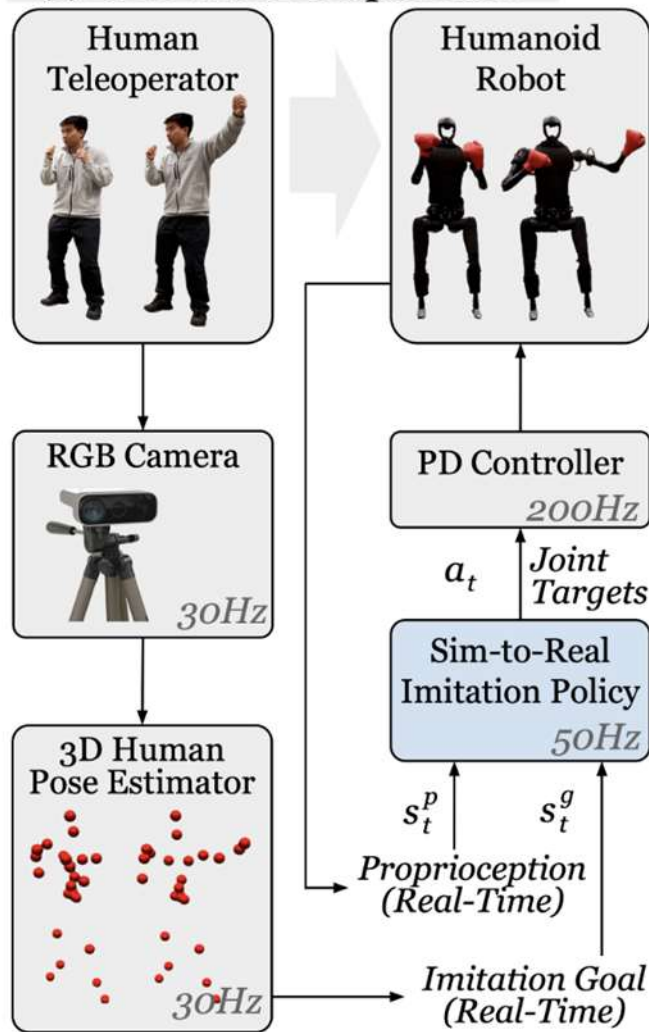
动作数据集
(干净且动作可行)

清除**不现实动作**后, 剩下的数据集会
被进一步清洗和优化。例如对动作序列
进行**平滑处理**以减少不连贯动作等。

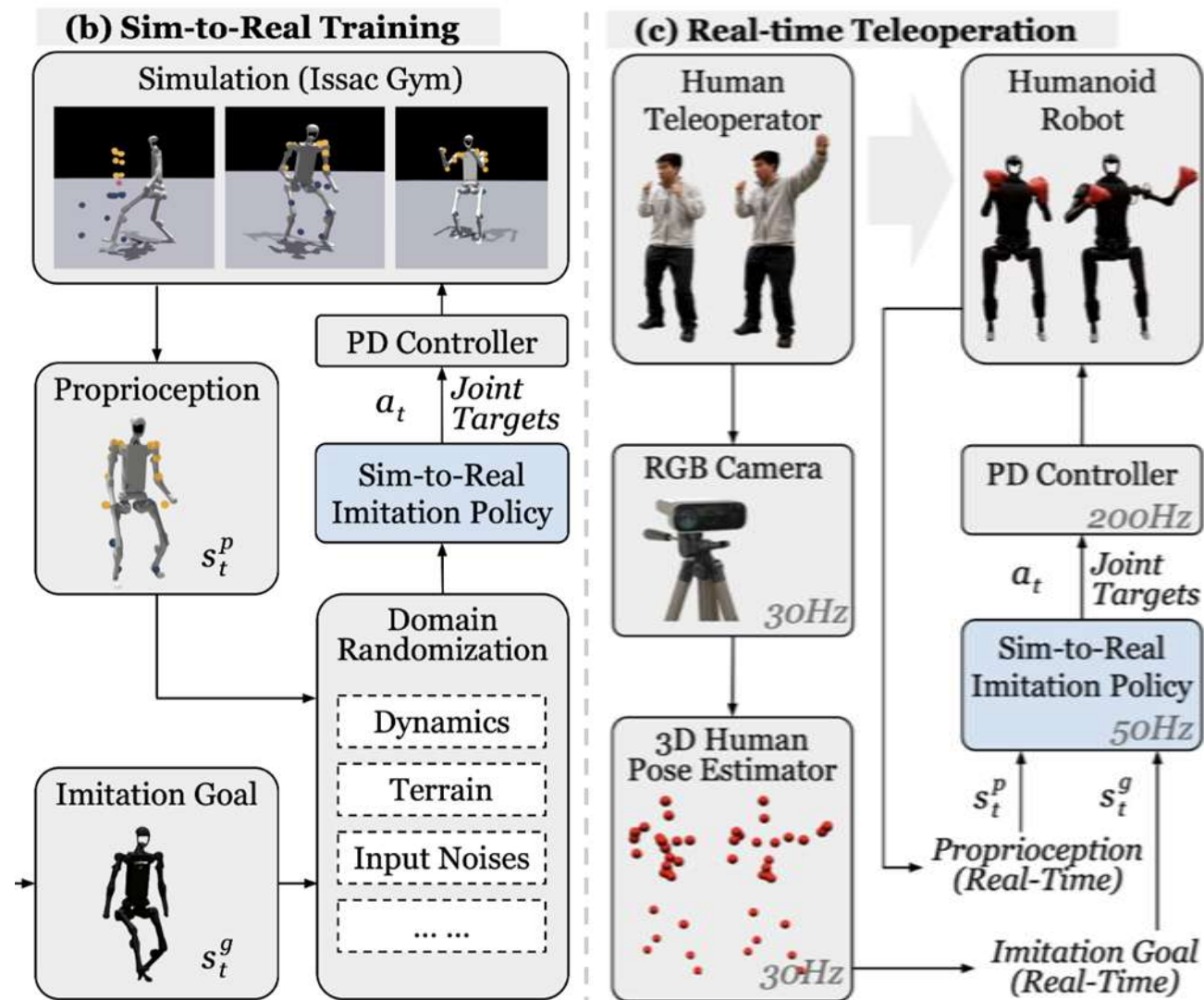
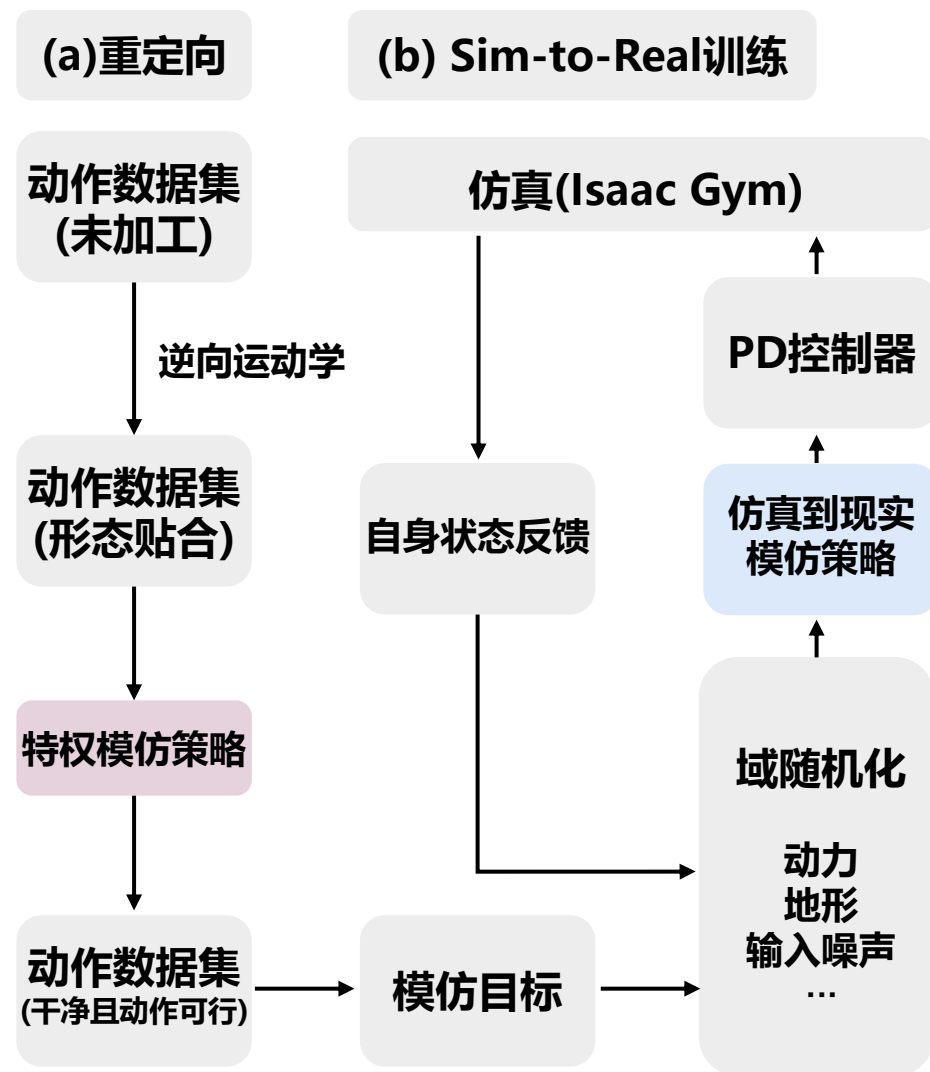
(b) Sim-to-Real Training



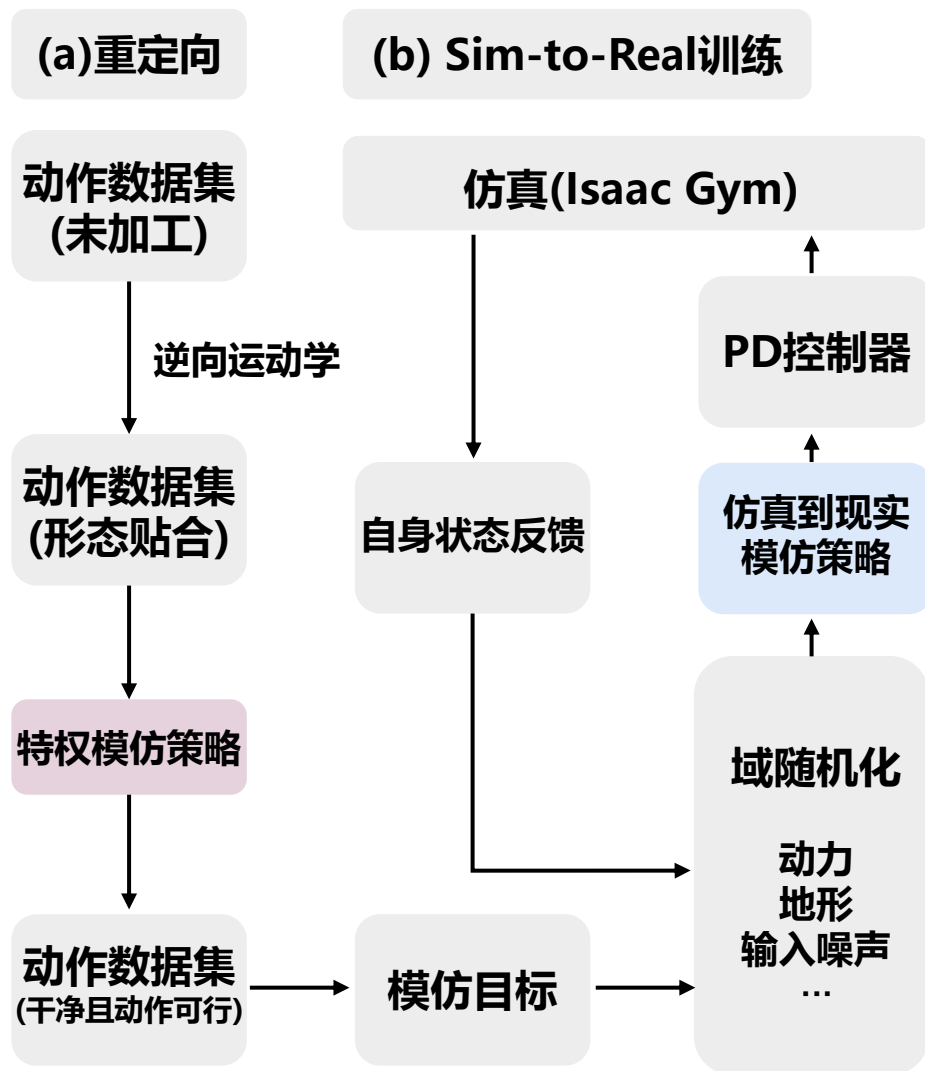
(c) Real-time Teleoperation



整体框架(包括abc三阶段)



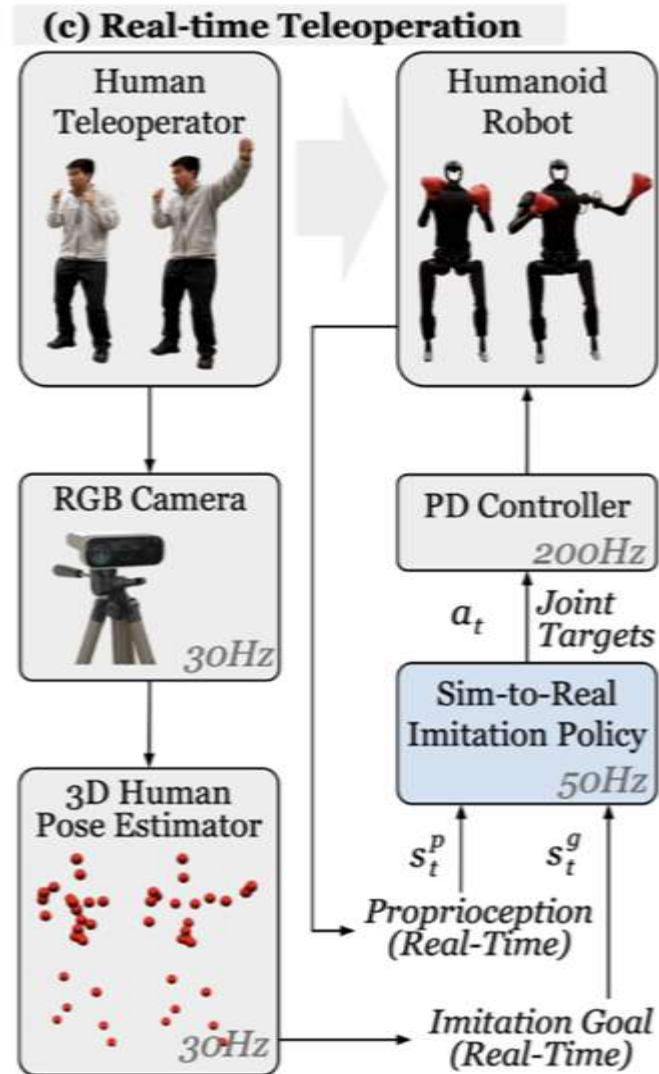
整体框架(包括abc三阶段)



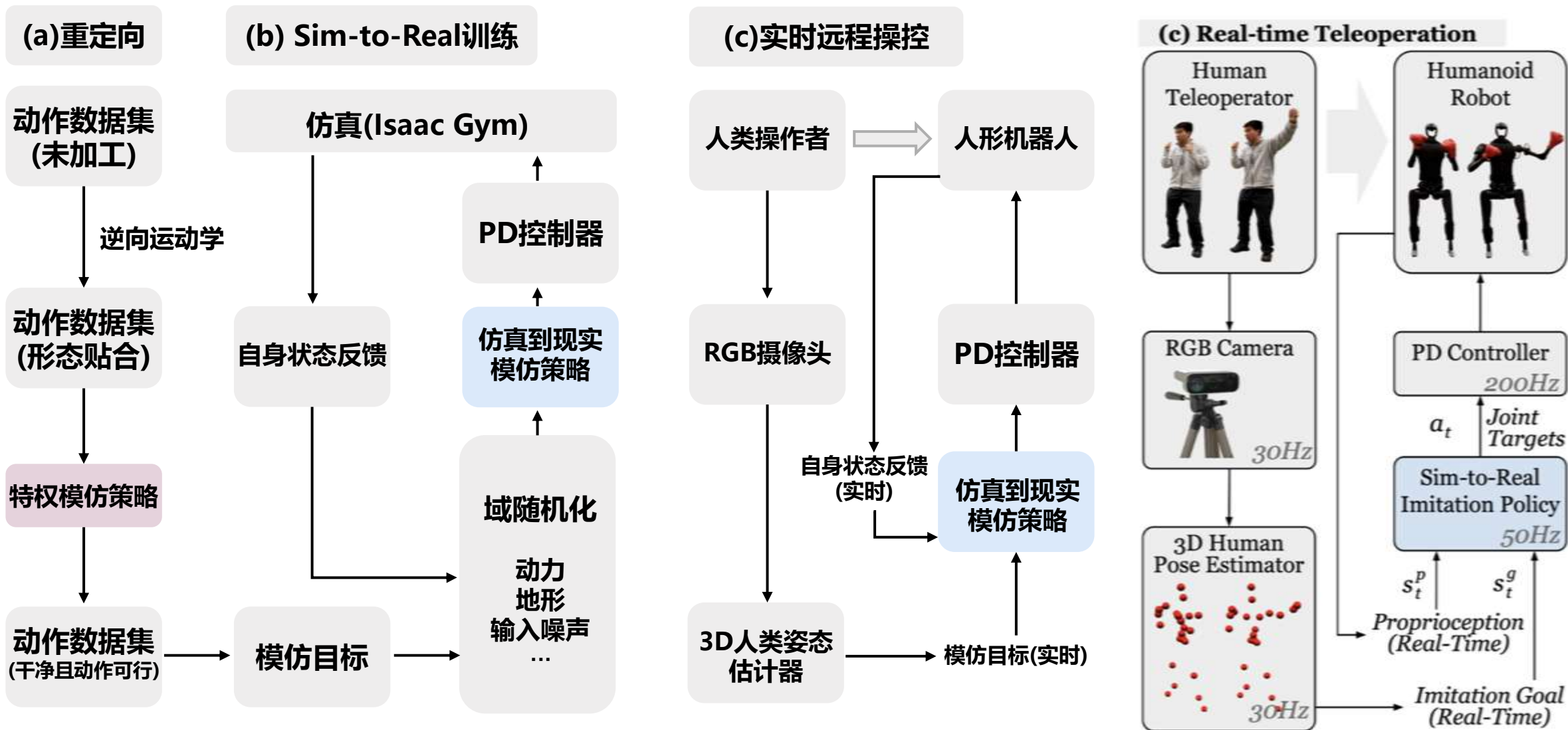
模仿策略的训练采用**强化学习**方法，会根据与人类动作的匹配程度来获得奖励。

策略会不断尝试新的行动，并根据结果(例如动作的准确性、机器人的稳定性等)来更新其参数，从而逐步提高模仿人类动作的能力。

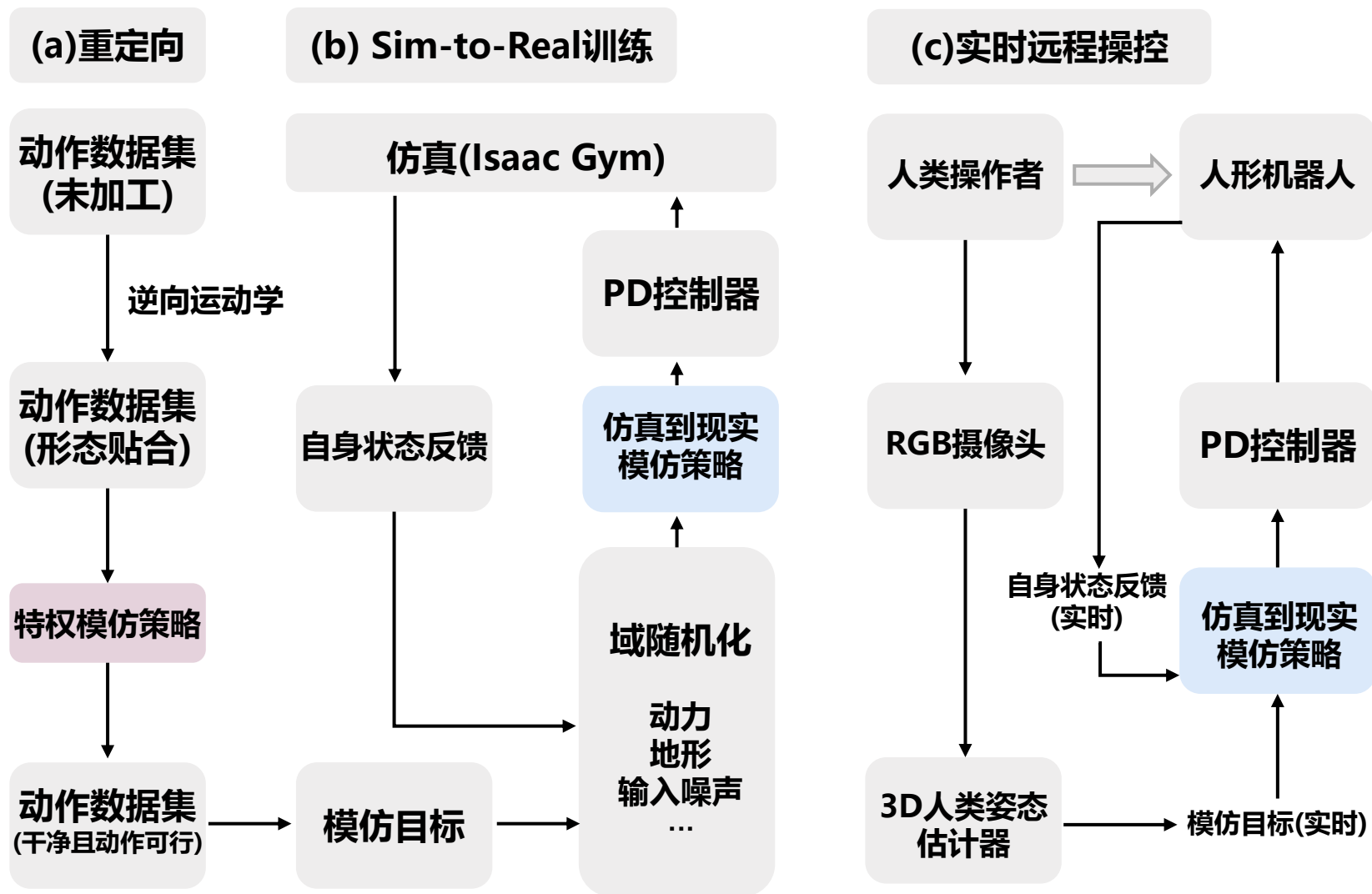
域随机化通过在仿真环境中随机改变某些参数(如摩擦系数、机器人质量、扭矩噪声等)来增加训练数据的多样性，以弥合仿真与现实之间的差距。



整体框架(包括abc三阶段)



整体框架(包括abc三阶段)



稳健性测试robustness

动作捕捉:

只需RGB摄像头进行动作捕捉, 姿态估计器能从摄像头捕获的视觉数据中提取出人体的关键点信息, 进而推断出操作者的整体姿态和运动。

2. HumanPlus-Stanford



系统特点

1. 实时模仿

允许使用单个RGB摄像头进行全身控制。

2. 自主任务

本质是模仿学习算法，能够通过40次人类演示，自主完成生活中的任务。如自主穿鞋、卸货、打字等。

RGB摄像头拍摄人类运动



将人类姿势重定向为机器人姿势

机器人本体感知到的信息

关节定位: 从SMPL-X到机器人

身体姿势估计: **WHAM**

手势估计: **HaMeR**

IMU和关节编码器测量关节角度、角速度等

目标姿态

当下姿态

输入

对机器人姿势进行条件训练(HST)

low-level策略

用PD控制器将目标角度转换为关节
扭矩, 控制机器人运动



2.HumanPlus-Stanford

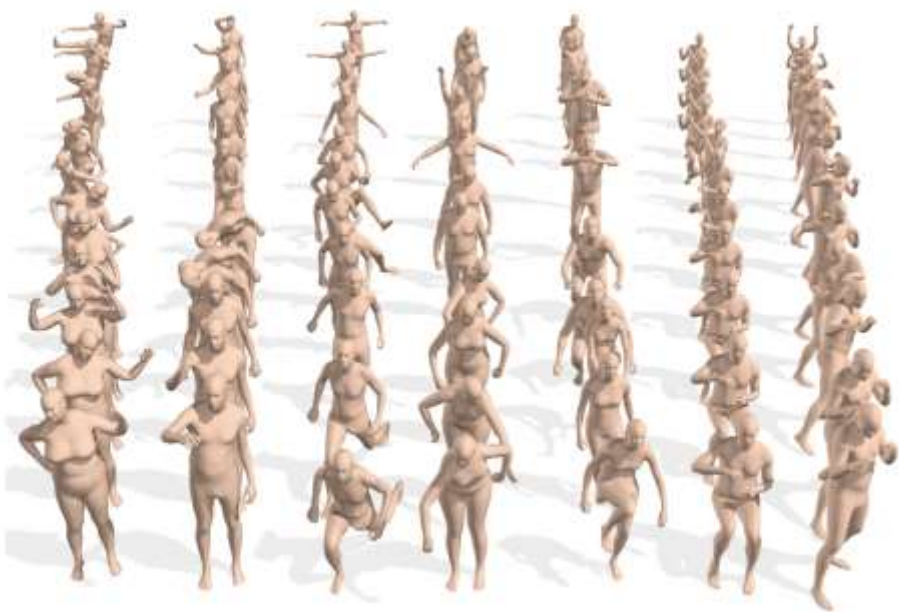


HumanPlus借鉴了H2O的Sim-to-Real的成功经验，使用大规模的强化学习训练用于全身控制的低级策略。

然而，基于学习的低级策略由于耗时的**奖励工程被设计为特定任务**，使得人形机器人一次只能展示一种技能，例如行走，这种限制限制了人形机器人平台能够执行的任务的多样性。

关键在于有一个40小时的**人体运动数据集AMASS**，涵盖了广泛的技能。使用该运动数据集，通过强化学习在模拟环境中(比如MuJoCo、Bullet等)训练一个低级策略。这个策略转移到现实世界，使人形机器人能够跟随人类身体和手部运动。

数据集的集成——AMASS



现有的人体动作捕捉数据集较小，动作有限，虽然有许多不同的数据集可用，但它们都对身体进行**不同参数化**，很难将它们汇集到单个数据集当中，在一项任务中使用。

AMASS是一个大型人体运动数据库，**统一**了 15 种不同的基于光学标记的动捕数据集。AMASS的一致表示使其可用于动画、可视化和生成深度学习的训练数据，拥有超过40小时的运动数据，涵盖300多个主题，超过11000个运动。

3. OpenTeleVision-UCSD



“相隔3000英里的遥操作”

远程操作系统

执行

使用关节映射来操纵机器人，要求操作员和机器人必须在同一地点，无法进行远程控制。

好在，VR头显通常集成内置手部跟踪算法，融合了来自多种传感器数据，包括多个摄像头、深度传感器和IMU。

通过VR设备收集的手部跟踪数据通常比自开发的视觉跟踪系统更稳定和准确，而后者仅使用了VR传感器的一部分(RGB+RGBD,Depth+IMU等)。

感知

对于感知，最直接的方法是操作员以第三人称或第一人称视角观察机器人任务空间。这不可避免地会在远程操作过程中被遮挡视线（如被机器人手臂或躯干遮挡）。

同时，对于精细的操作任务，远程操作员难以近距离直观地观察物体。

创新点：单一设备同时做到既可远程控制又可深度感知

主动视觉反馈



第一人称主动感知：机器人头部具备主动立体RGB相机，配备2或3个自由度的驱动，摄像头会随着操作员的头部移动而移动，进行流媒体传输。

机械臂、机械手控制



手臂控制：基于Pinocchio的闭环逆运动学算法计算关节角
手部控制：通过运动重定向库dex-retargeting，将人手关键点转换为机器人关节角度命令

3. OpenTeleVision-UCSD

远程操作

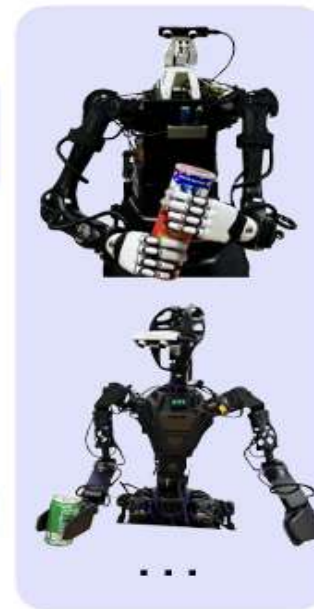


立体视觉

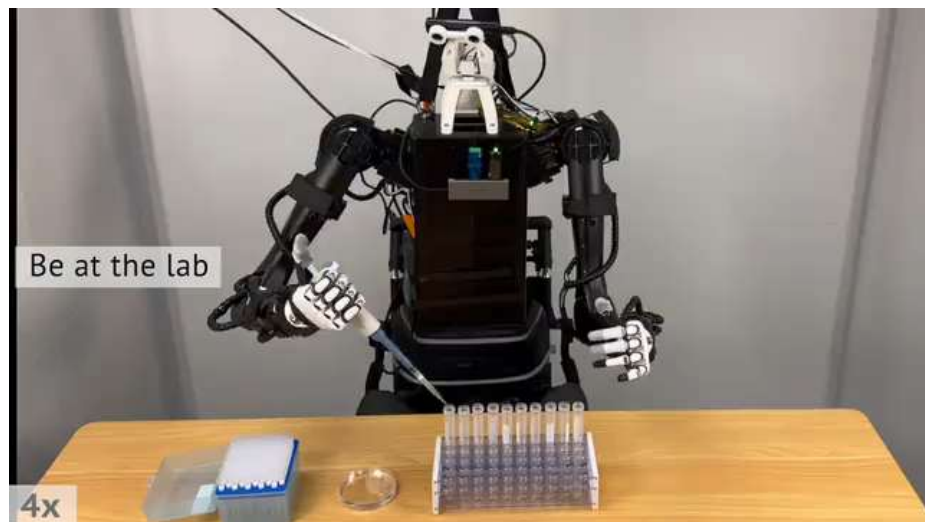


60Hz

手，头，腕姿势



遥操作下的精细工作



移液管任务表明OpenTeleVision也能够执行精确的动作。这也是一个不可能由机械手夹持器完成的任务，因为移液管的使用是专为类人手设计的。

多平台兼容的系统



在VR设备上，用户可以通过手部和手腕姿态流远程操作。在其他设备上，手和手腕的流媒体传输不可用，但用户可以通过在设备屏幕上拖动来控制机器人的主动颈部，看到传输的图像。

作者	设计	图片	总结
Tairan He等 “H2O” 2024.3	1.重定向→ Sim-to-Real→实时遥操作 2.仅用RGB摄像头实现动作捕捉 3.40小时人体运动数据集AMASS 4.Sim-to-Real训练设计		一个基于强化学习的系统，可实现仅使用 RGB相机的人形机器人的实时全身远程操作
Zipeng Fu等 “HumanPlus” 2024.6	1.重定向→条件训练低级策略→实时遥操作 2.仅用RGB摄像头实现动作捕捉 3.40小时人体运动数据集AMASS 4.HST条件训练设计		和H2O的硬件，逻辑基本相同，采用另外一种强化学习的策略。能学习人类演示自主完成任务
Xuxin Cheng等 “OpenTeleVision” 2024.7	1.操作者穿戴VR头显，机器人头部具备主动立体相机 2.头部可跟随运动，主动感知 3.远程控制：操作者不必在现场 4.深度感知：避免第三人称视角的弊端，实现细粒度操作		真实虚拟体验、远程现实控制兼备。感知端人类第一人称视角沉浸式体验；执行端突破距离限制，真正意义上“遥”操作

[1]He T, Luo Z, Xiao W, et al. Learning human-to-humanoid real-time whole-body teleoperation[J]. arXiv preprint arXiv:2403.04436, 2024.
[2]Fu Z, Zhao Q, Wu Q, et al. HumanPlus: Humanoid Shadowing and Imitation from Humans[J]. arXiv preprint arXiv:2406.10454, 2024.
[3]Cheng X, Li J, Yang S, et al. Open-TeleVision: teleoperation with immersive active visual feedback[J]. arXiv preprint arXiv:2407.01512, 2024.