Multilevel models are extensions of standard models designed to deal with the problem of omitted variable bias. No statistical model of a social system is ever complete because the system is complex and highly interactive. People live in and interact with their environments in complex ways that will often distort traditional regression type model. In the weakest case this presents as an i.i.d (*independent and identically distributed error*) violation that can often produce incorrect standard errors that are usually, but not always, too small. In the broadest case it can cause you to misunderstand substantive relationships between Y and X because effects over subpopulations are inappropriately averaged together. This is fundamentally no different than any other important omitted variable bias problem or missing data problem.

The problem practical researchers have to deal with is how to estimate a model that takes into account heterogeneity across people/places/actors/firms without knowing what is causing it in the first place. In panel data econometrics we know that responses across time for a single person are likely to be highly correlated. There are entire academic disciplines devoted to understanding the underlying causes and it is safe to say that no single statistical model is going to adequately explain all of the within person correlation as a function of characteristics of people. So we have to use the people themselves as variables. We can do the same thing with context by accounting for social networks or geography.

**Fixed vs. Random Effects**

*Econometric fixed effects*: A collection of subject/group specific intercepts by way of a set of indicator variables. You can see the equation here: $Y_{ij} = \alpha + \beta X_{ij} + \beta Z_j + \varepsilon$. Y is some outcome indexed by both individual observations i and groups j. $\alpha$ is a grand mean or intercept. The X's are also indexed by individuals and groups and the Z's are indicator variables for each group minus a reference. This model removes between group variation by directly including the indicator variables in the model so that the only thing left for the X's to explain is within group variability.

Econometric fixed effects are useful in a number of applied settings:
- If you are interested in the values of $\beta Z_j$ directly
  - To test differences in groups like you would with an Anova
  - To rank order groups based on their expected value of Y conditional on other things
- If $\beta Z_j$ represents some actual mistake in either your research design or the implementation of your design. Here, between group effects are pathological to your research design and cutting them out of your data helps you make better inferences.
  - If you are doing a controlled experiment with some form of random assignment but the randomization process did not magically lead to equivalent groups—it usually doesn't.
  - If you are studying some kind of observational treatment effect and need to pretend you have something like random assignment. This is why fixed effects estimators are often classified under tools for causal inference.
- If you are interested in potentially mixing heterogeneous populations within an analysis but do not have enough groups to effectively use random effects
  - e.g. using dummy variables for gender or race instead of attempting incredibly ill-fitting random effects

Econometric fixed effects are not useful when:
- You are actually trying to model the processes or mechanisms that cause Y based on observational data. In that case fixed effects will actually obscure the effects of any between group variables on Y as well as their interactions with other variables. This can and usually does lead to a fundamental misunderstanding about the factors that drive an outcome variable.

*Econometric random effects*: A single (usually latent) variable that is estimated as the expected value of Y for each group. Random effects are literally made from fixed effects. They are not different conceptual objects. They are the same information with a different functional form. Using fixed effects or random effects is a research design question and not fundamentally a statistical one.[1]

People often think that it IS a statistical question as to if there is a fixed or random effect and that the Hausman test (Amini et al. 2012; Nerlove 2000) can tell them the answer. The Hausman test compares a set of coefficients from a fixed effects model and a random effects model and tells you which coefficients changed. If there is any change, people often mistakenly believe it means the FE estimator is appropriate. If the FE estimator was appropriate to what you wanted to study then there would never be any need to run a Hausman test in the first place because the RE estimator couldn't give you the appropriate information to begin with. What it actually means is that between group variance captured by $\beta Z_j$ (aka the dummy variables) is related to variance in $\beta X_{ij}$ (aka your variables of interest) and so you need to use the Mundlak (1978) specification to explicitly model it. The Hausman test is telling you that you have omitted time or group invariant variables that are related to your X's and thus omitted variable bias.

The Hausman test actually means that $X_{ij} \perp Z_j$. If that is not the case then you need to use a random effects model with group level averages of X or $\bar{X}_j$ in the model and transform $X_{ij}$ by subtracting $\bar{X}_j$ from it. This effectively means you pass the Hausman test because your $X_i's$ are now uncorrelated with any group level variables including the random effect. At that point you get $Y_{ij} = \alpha + \beta(X_{ij} - \bar{X}_j) + \beta\bar{X}_j + \mu + \varepsilon$. In this model $(X_{ij} - \bar{X}_j)$ is functionally equivalent to the $X_i$ of an econometric fixed effects estimator. The $\bar{X}_j$ term means group averages (though in very special some cases you could potentially use medians or some combination of splines or polynomials). The $\mu$ term signifies a random effect or varying intercept around $\alpha$ that signifies group-specific intercepts equivalent to the fixed effects dummy variable coefficients in a FE model. The $\varepsilon$ is the non-group related residual error term.

*Statistical fixed effects*: any estimated component of a model that does not vary over subgroups. In a traditional linear or nonlinear model (like OLS or logit) all effects are fixed meaning there is only one coefficient for alpha and each beta. A random intercept model changes this by allowing a range of group-specific intercepts and thus no longer fixing a universal intercept.

*Statistical random effects*: any component of a model that varies over subgroups. This means a random intercept but it also means random coefficients which are interactions between the random intercept and a within-group varying variable. So the statistical literature can refer to both a random effects model and a random effect within a model.

---

[1] With the exception of using fixed effects to control for heterogeneous subpopulations when you don't have enough groups to use a random effect.